

Springer Texts in Statistics

Kenneth Lange

Optimization

Second Edition



Springer

Kenneth Lange

Optimization

Second Edition

 Springer

Kenneth Lange
Biomathematics, Human Genetics,
Statistics
University of California
Los Angeles, CA, USA

STS Editorial Board

George Casella
Department of Statistics
University of Florida
Gainesville, FL, USA

Ingram Olkin
Department of Statistics
Stanford University
Stanford, CA, USA

Stephen Fienberg
Department of Statistics
Carnegie Mellon University
Pittsburg, PA, USA

ISSN 1431-875X
ISBN 978-1-4614-5837-1 ISBN 978-1-4614-5838-8 (eBook)
DOI 10.1007/978-1-4614-5838-8
Springer New York Heidelberg Dordrecht London

Library of Congress Control Number: 2012948598

© Springer Science+Business Media New York 2004, 2013

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed. Exempted from this legal reservation are brief excerpts in connection with reviews or scholarly analysis or material supplied specifically for the purpose of being entered and executed on a computer system, for exclusive use by the purchaser of the work. Duplication of this publication or parts thereof is permitted only under the provisions of the Copyright Law of the Publisher's location, in its current version, and permission for use must always be obtained from Springer. Permissions for use may be obtained through RightsLink at the Copyright Clearance Center. Violations are liable to prosecution under the respective Copyright Law.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

While the advice and information in this book are believed to be true and accurate at the date of publication, neither the authors nor the editors nor the publisher can accept any legal responsibility for any errors or omissions that may be made. The publisher makes no warranty, express or implied, with respect to the material contained herein.

Printed on acid-free paper

Springer is part of Springer Science+Business Media (www.springer.com)

To my daughters, Jane and Maggie

Preface to the Second Edition

If this book were my child, I would say that it has grown up after passing through a painful adolescence. It has rebelled at times over my indecision. While it has aspired to greatness, I have aspired to closure. Although the second edition of *Optimization* occasionally demanded more exertion than I could muster, writing it has broadened my intellectual horizons. I hope my own struggles to reach clarity will translate into an easier path for readers.

The book's stress on mathematical fundamentals continues. Indeed, I was tempted to re-title it *Mathematical Analysis and Optimization*. I resisted this temptation because my ultimate goal is still to teach optimization theory. Nonetheless, there is a new chapter on the gauge integral and expanded treatments of differentiation and convexity. The focus remains on finite-dimensional optimization. The sole exception to this rule occurs in the new chapter on the calculus of variations. In my view functional analysis is just too high a rung on the ladder of mathematical abstraction.

Covering all of optimization theory is simply out of the question. Even though the second edition is more than double the length of the first, many important topics are omitted. The most grievous omissions are the simplex algorithm of linear programming and modern interior point methods. Fortunately, there are many admirable books devoted to these subjects. My development of adaptive barrier methods and exact penalty methods also partially compensates.

In addition to the two chapters on integration and the calculus of variations, four new chapters treat block relaxation (block descent and block ascent) and various advanced topics in the convex calculus, including the

Fenchel conjugate, subdifferentials, duality, feasibility, alternating projections, projected gradient methods, exact penalty methods, and Bregman iteration. My own interests in data mining and biological applications have dictated the nature of these chapters. High-dimensional problems are driving the discipline of optimization. These are qualitatively different from traditional problems, and standard algorithms such as Newton's method are often impractical. Penalization, model sparsity, and the MM algorithm now assume dominant roles. Fortunately, many of the challenging modern problems can also be phrased as convex programs.

In the first edition I eschewed the convention of setting vectors and matrices in boldface type. In the second edition I embrace it. Although this decision improves readability, it carries with it some residual ambiguity. The main difficulty lies in distinguishing constant vectors and matrices from vector and matrix-valued functions. In general, I have elected to set functions in ordinary type even when they are vector or matrix valued. The exceptions occur in the calculus of variations, where functions are considered vectors in infinite-dimensional spaces. Thus, a function appears in ordinary type when its argument is displayed and in boldface type when its argument is omitted.

Many people have helped me prepare this second edition. Hua Zhou and Tongtong Wu, my former postdoctoral fellows, and Eric Chi, my current postdoctoral fellow, deserve special credit. Without their assistance, the book would have been intellectually duller and graphically drearier. I would also like to thank my former doctoral students David Alexander, David Hunter, Mary Sehl, and Jinjin Zhou for proofreading and critiquing the new material. The students in my optimization class checked most of the exercises. I am indebted to Forrest Crawford, Gabriela Cybis, Gary Evans, Mitchell Johnson, Wesley Kerr, Kevin Keys, Omid Kohannim, Lewis Lee, Matthew Levinson, Lae Un Kim, John Ranola, and Moses Wilkes for their help.

Finally, let me report on my daughters Maggie and Jane, to whom this book is dedicated. Maggie is now embarked on a postdoctoral fellowship in medical ethics at Macquarie University in Sydney, Australia. Jane is completing her dissertation in biostatistics at the University of Washington. Assimilating their scholarship will keep me young for many years to come.

Preface to the First Edition

This foreword, like many forewords, was written afterwards. That is just as well because the plot of the book changed during its creation. It is painful to recall how many times classroom realities forced me to shred sections and start anew. Perhaps such adjustments are inevitable. Certainly I gained a better perspective on the subject over time. I also set out to teach optimization theory and wound up teaching mathematical analysis. The students in my classes are no less bright and eager to learn about optimization than they were a generation ago, but they tend to be less prepared mathematically. So what you see before you is a compromise between a broad survey of optimization theory and a textbook of analysis. In retrospect, this compromise is not so bad. It compelled me to revisit the foundations of analysis, particularly differentiation, and to get right to the point in optimization theory.

The content of courses on optimization theory varies tremendously. Some courses are devoted to linear programming, some to nonlinear programming, some to algorithms, some to computational statistics, and some to mathematical topics such as convexity. In contrast to their gaps in mathematics, most students now come well trained in computing. For this reason, there is less need to emphasize the translation of algorithms into computer code. This does not diminish the importance of algorithms, but it does suggest putting more stress on their motivation and theoretical properties. Fortunately, the dichotomy between linear and nonlinear programming is fading. It makes better sense pedagogically to view linear programming as a special case of nonlinear programming. This is the attitude taken in

the current book, which makes little mention of the simplex method and develops interior point methods instead. The real bridge between linear and nonlinear programming is convexity. I stress not only the theoretical side of convexity but also its applications in the design of algorithms for problems with either large numbers of parameters or nonlinear constraints.

This graduate-level textbook presupposes knowledge of calculus and linear algebra. I develop quite a bit of mathematical analysis from scratch and feature a variety of examples from linear algebra, differential equations, and convexity theory. Of course, the greater the prior exposure of students to this background material, the more quickly the beginning chapters can be covered. If the need arises, I recommend the texts [82, 134, 135, 188, 222, 223] for supplementary reading. There is ample material here for a fast-paced, semester-long course. Instructors should exercise their own discretion in skipping sections or chapters. For example, Chap. 10 on the EM algorithm primarily serves the needs of students in biostatistics and statistics. Overall, my intended audience includes graduate students in applied mathematics, biostatistics, computational biology, computer science, economics, physics, and statistics. To this list I would like to add upper-division majors in mathematics who want to see some rigorous mathematics with real applications. My own background in computational biology and statistics has obviously dictated many of the examples in the book.

Chapter 1 starts with a review of exact methods for solving optimization problems. These are methods that many students will have seen in calculus, but repeating classical techniques with fresh examples tends simultaneously to entertain, instruct, and persuade. Some of the exact solutions also appear later in the book as parts of more complicated algorithms.

Chapters 2 through 4 review undergraduate mathematical analysis. Although much of this material is standard, the examples may keep the interest of even the best students. Instructors should note that Carathéodory's definition rather than Fréchet's definition of differentiability is adopted. This choice eases the proof of many results. The gauge integral, another good addition to the calculus curriculum, is mentioned briefly.

Chapter 5 gets down to the serious business of optimization theory. McShane's clever proof of the necessity of the Karush–Kuhn–Tucker conditions avoids the complicated machinery of manifold theory and convex cones. It makes immediate use of the Mangasarian–Fromovitz constraint qualification. To derive sufficient conditions for optimality, I introduce second differentials by extending Carathéodory's definition of first differentials. To my knowledge, this approach to second differentials is new. Because it melds so effectively with second-order Taylor expansions, it renders critical proofs more transparent.

Chapter 6 treats convex sets, convex functions, and the relationship between convexity and the multiplier rule. The chapter concludes with the derivation of some of the classical inequalities of probability theory. Prior exposure to probability theory will obviously be an asset for readers here.

Chapters 8 and 9 introduce the MM and EM algorithms. These exploit convexity and the notion of majorization in transferring minimization of the objective function to a surrogate function. Minimizing the surrogate function drives the objective function downhill. The EM algorithm, which is a special case of the MM algorithm, arose in statistics. It is a slight misnomer to call these algorithms. They are really prescriptions for constructing algorithms. It takes experience and skill to wield these tools effectively, so careful attention to the examples is imperative.

Chapter 10 covers Newton's method and its statistical variants, scoring and the Gauss–Newton algorithm. To make this material less dependent on statistical knowledge, I have tried to motivate several algorithms from the perspective of positive definite approximation of the second differential of the objective function. Chapter 11 covers the conjugate gradient algorithm, quasi-Newton algorithms, and the method of trust regions. These classical subjects are in danger of being dropped from the curriculum of nonlinear programming. In my view, this would be a mistake.

Chapter 12 is devoted to convergence questions, both local and global. This material beautifully illustrates the virtues of soft analysis. Instructors wanting to emphasize practical matters may be tempted to sacrifice Chap. 12, but the constant interplay between theory and practice in designing new algorithms argues for its inclusion.

Chapter 13 on convex programming ends the book where more advanced treatises would start. I discuss adaptive barrier methods as a novel application of the MM algorithm, Dykstra's algorithm for finding feasible points in convex programming, and the rudiments of duality theory. These topics belong to the promised land. All you get here is a glimpse from the mountaintop looking out across the river.

Let me add a few words about notation. Lower-division undergraduate texts carefully distinguish between scalars and vectors by setting vectors in boldface type. This convention is considered cumbersome in higher mathematics and is dropped. However, mathematical analysis is plagued by a proliferation of superscripts and subscripts. I prefer to avoid superscripts because of the possible confusion with powers. This decision makes it difficult to distinguish an element of a vector sequence from a component of a vector. My compromise is to represent the m th entry of a vector sequence as $x_{(m)}$ and the n th component of that sequence element as x_{mn} . Similar conventions hold for matrices. Thus, M_{jkl} is the entry in row k and column l of the j th matrix $M_{(j)}$ of a sequence of matrices. Elements of scalar sequences are subscripted in the usual fashion without the enclosing parentheses.

I would like to thank my UCLA students for their help and patience in debugging this text. If it is readable, it is because their questions cut through the confusion. In retrospect, there were more contributing students than I can credit. Let me single out Jason Aten, Lara Bauman, Brian Dolan,

Wei-Hsun Liao, Andrew Nevai-Tucker, Robert Rovetti, and Andy Yip. Paul Maranian kindly prepared the index and proofread my last draft. Finally, I thank my ever helpful and considerate editor, John Kimmel.

I dedicate this book to my daughters, Jane and Maggie. It has been a privilege to be your father. Now that you are adults, I hope you can find the same pleasure in pursuing ideas that I have found in my professional life.

Contents

Preface to the Second Edition	vii
Preface to the First Edition	ix
1 Elementary Optimization	1
1.1 Introduction	1
1.2 Univariate Optimization	1
1.3 Multivariate Optimization	7
1.4 Constrained Optimization	10
1.5 Problems	17
2 The Seven C's of Analysis	23
2.1 Introduction	23
2.2 Vector and Matrix Norms	23
2.3 Convergence and Completeness	26
2.4 The Topology of \mathbb{R}^n	30
2.5 Continuous Functions	34
2.6 Semicontinuity	42
2.7 Connectedness	44
2.8 Uniform Convergence	46
2.9 Problems	47
3 The Gauge Integral	53
3.1 Introduction	53
3.2 Gauge Functions and δ -Fine Partitions	54
	xiii

3.3	Definition and Basic Properties of the Integral	57
3.4	The Fundamental Theorem of Calculus	62
3.5	More Advanced Topics in Integration	66
3.6	Problems	71
4	Differentiation	75
4.1	Introduction	75
4.2	Univariate Derivatives	75
4.3	Partial Derivatives	79
4.4	Differentials	81
4.5	Multivariate Mean Value Theorem	88
4.6	Inverse and Implicit Function Theorems	89
4.7	Differentials of Matrix-Valued Functions	93
4.8	Problems	98
5	Karush-Kuhn-Tucker Theory	107
5.1	Introduction	107
5.2	The Multiplier Rule	108
5.3	Constraint Qualification	114
5.4	Taylor-Made Higher-Order Differentials	117
5.5	Applications of Second Differentials	123
5.6	Problems	128
6	Convexity	137
6.1	Introduction	137
6.2	Convex Sets	138
6.3	Convex Functions	142
6.4	Continuity, Differentiability, and Integrability	149
6.5	Minimization of Convex Functions	152
6.6	Moment Inequalities	159
6.7	Problems	162
7	Block Relaxation	171
7.1	Introduction	171
7.2	Examples of Block Relaxation	172
7.3	Problems	180
8	The MM Algorithm	185
8.1	Introduction	185
8.2	Philosophy of the MM Algorithm	186
8.3	Majorization and Minorization	187
8.4	Allele Frequency Estimation	189
8.5	Linear Regression	191
8.6	Bradley-Terry Model of Ranking	193
8.7	Linear Logistic Regression	194

8.8	Geometric and Signomial Programs	194
8.9	Poisson Processes	197
8.10	Transmission Tomography	198
8.11	Poisson Multigraphs	202
8.12	Problems	204
9	The EM Algorithm	221
9.1	Introduction	221
9.2	Definition of the EM Algorithm	222
9.3	Missing Data in the Ordinary Sense	224
9.4	Allele Frequency Estimation	225
9.5	Clustering by EM	226
9.6	Transmission Tomography	228
9.7	Factor Analysis	230
9.8	Hidden Markov Chains	234
9.9	Problems	236
10	Newton’s Method and Scoring	245
10.1	Introduction	245
10.2	Newton’s Method and Root Finding	246
10.3	Newton’s Method and Optimization	248
10.4	MM Gradient Algorithm	250
10.5	Ad Hoc Approximations of $d^2 f(\theta)$	252
10.6	Scoring and Exponential Families	254
10.7	The Gauss-Newton Algorithm	257
10.8	Generalized Linear Models	258
10.9	Accelerated MM	259
10.10	Problems	262
11	Conjugate Gradient and Quasi-Newton	273
11.1	Introduction	273
11.2	Centers of Spheres and Centers of Ellipsoids	274
11.3	The Conjugate Gradient Algorithm	275
11.4	Line Search Methods	278
11.5	Stopping Criteria	280
11.6	Quasi-Newton Methods	281
11.7	Trust Regions	285
11.8	Problems	286
12	Analysis of Convergence	291
12.1	Introduction	291
12.2	Local Convergence	292
12.3	Coercive Functions	297
12.4	Global Convergence of the MM Algorithm	299
12.5	Global Convergence of Block Relaxation	302

12.6	Global Convergence of Gradient Algorithms	303
12.7	Problems	306
13	Penalty and Barrier Methods	313
13.1	Introduction	313
13.2	Rudiments of Barrier and Penalty Methods	314
13.3	An Adaptive Barrier Method	318
13.4	Imposition of a Prior in EM Clustering	325
13.5	Model Selection and the Lasso	327
13.6	Lasso Penalized ℓ_1 Regression	329
13.7	Lasso Penalized ℓ_2 Regression	330
13.8	Penalized Discriminant Analysis	333
13.9	Problems	334
14	Convex Calculus	341
14.1	Introduction	341
14.2	Notation	342
14.3	Fenchel Conjugates	342
14.4	Subdifferentials	351
14.5	The Rules of Convex Differentiation	358
14.6	Spectral Functions	365
14.7	A Convex Lagrange Multiplier Rule	372
14.8	Problems	375
15	Feasibility and Duality	383
15.1	Introduction	383
15.2	Dykstra's Algorithm	384
15.3	Contractive Maps	389
15.4	Dual Functions	393
15.5	Examples of Dual Programs	396
15.6	Practical Applications of Duality	402
15.7	Problems	406
16	Convex Minimization Algorithms	415
16.1	Introduction	415
16.2	Projected Gradient Algorithm	416
16.3	Exact Penalties and Lagrangians	421
16.4	Mechanics of Path Following	426
16.5	Bregman Iteration	432
16.6	Split Bregman Iteration	436
16.7	Convergence of Bregman Iteration	439
16.8	Problems	440
17	The Calculus of Variations	445
17.1	Introduction	445
17.2	Classical Problems	446

17.3	Normed Vector Spaces	448
17.4	Linear Operators and Functionals	451
17.5	Differentials	453
17.6	The Euler-Lagrange Equation	456
17.7	Applications of the Euler-Lagrange Equation	459
17.8	Lagrange's Lacuna	462
17.9	Variational Problems with Constraints	464
17.10	Natural Cubic Splines	466
17.11	Problems	467
Appendix: Mathematical Notes		473
A.1	Univariate Normal Random Variables	473
A.2	Multivariate Normal Random Vectors	475
A.3	Polyhedral Sets	477
A.4	Birkhoff's Theorem and Fan's Inequality	480
A.5	Singular Value Decomposition	485
A.6	Hadamard Semidifferentials	487
A.7	Problems	497
References		499
Index		519

1

Elementary Optimization

1.1 Introduction

As one of the oldest branches of mathematics, optimization theory served as a catalyst for the development of geometry and differential calculus [258]. Today it finds applications in a myriad of scientific and engineering disciplines. The current chapter briefly surveys material that most students encounter in a good calculus course. This review is intended to showcase the variety of methods used to find the exact solutions of elementary problems. We will return to some of these methods later from a more rigorous perspective. One of the recurring themes in optimization theory is its close connection to inequalities. This chapter introduces a few classical inequalities; more will appear in succeeding chapters.

1.2 Univariate Optimization

The first optimization problems students encounter are univariate. Solution techniques for these simple problems are hardly limited to differential calculus. Our first two examples illustrate how plane geometry and algebra can play a role.

Example 1.2.1 *Heron's Problem*

The ancient mathematician Heron of Alexandria posed one of the earliest optimization problems. Consider the two points A and B and the line

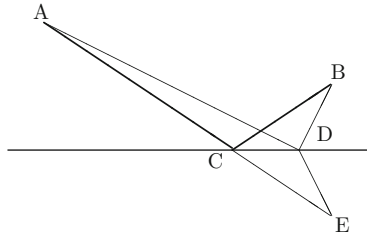


FIGURE 1.1. Diagram for Heron's problem

containing the points C and D drawn in Fig. 1.1. Heron's problem is to find the position of C on the line that minimizes the sum of the distances $|AC|$ and $|BC|$. The correct choice of C is determined by reflecting B across the given line to give E . From E we draw the line to A and note its intersection C with the original line. To demonstrate that C minimizes the total distance $|AC| + |BC|$, consider any other point D on the original horizontal line. By symmetry, $|AC| + |BC| = |AC| + |CE|$. Similarly, by symmetry, $|AD| + |BD| = |AD| + |DE|$. Because the sum of the lengths of two sides of a triangle exceeds the length of the third side, it follows immediately that $|AC| + |CE| \leq |AD| + |DE|$. Thus, C solves Heron's problem.

This example also has an optical interpretation. If we imagine that the horizontal line containing C lies on a mirror, then light travels along the quickest path between A and B via the mirror. This extremal principle can be explained by considering the wave nature of light, but we omit the long digression. It is interesting that the geometric argument automatically implies that the angle of incidence of the light ray equals the angle of reflection. ■

Example 1.2.2 Simple Arithmetic-Geometric Mean Inequality

If x and y are two nonnegative numbers, then $\sqrt{xy} \leq (x + y)/2$. This can be proved by noting that

$$\begin{aligned} 0 &\leq (\sqrt{x} - \sqrt{y})^2 \\ &= x - 2\sqrt{xy} + y. \end{aligned}$$

Evidently, equality holds if and only if $x = y$. As an application consider maximization of the function $f(x) = x(1 - x)$. The inequality just derived shows that

$$f(x) \leq \left(\frac{x + 1 - x}{2} \right)^2 = \frac{1}{4},$$

with equality when $x = 1/2$. Thus, the maximum of $f(x)$ occurs at the point $x = 1/2$. One can interpret $f(x)$ as the area of a rectangle of fixed

perimeter 2 with sides of length x and $1 - x$. The rectangle with the largest area is a square. The function $2f(x)$ is interpreted in population genetics as the fraction of a population that is heterozygous at a genetic locus with two alleles having frequencies x and $1 - x$. Heterozygosity is maximized when the two alleles are equally frequent. ■

With the advent of differential calculus, it became possible to solve optimization problems more systematically. Before discussing concrete examples, it is helpful to review some of the standard theory. We restrict attention to real-valued functions defined on intervals. The intervals in question can be finite or infinite in extent and open or closed at either end. According to a celebrated theorem of Weierstrass, a continuous function $f(x)$ defined on a closed finite interval $[a, b]$ attains its minimum and maximum values on the interval. These extremal values are necessarily finite. The extremal points can occur at the endpoints a or b or at an interior point c . In the latter case, when $f(x)$ is differentiable, an even older principle of Fermat requires that $f'(c) = 0$. The stationarity condition $f'(c) = 0$ is no guarantee that c is optimal. It is possible for c to be a local rather than a global minimum or maximum or even to be a saddle point. However, it usually is a simple matter to check the endpoints a and b and any stationary points c . Collectively, these points are known as critical points.

If the domain of $f(x)$ is not a closed finite interval $[a, b]$, then the minimum or maximum of $f(x)$ may not exist. One can usually rule out such behavior by examining the limit of $f(x)$ as x approaches an open boundary. For example on the interval $[a, \infty)$, if $\lim_{x \rightarrow \infty} f(x) = \infty$, then we can be sure that $f(x)$ possesses a minimum on the interval, and we can find it by comparing the values of $f(x)$ at a and any stationary points c . On a half open interval such as $(a, b]$, we can likewise find a minimum whenever $\lim_{x \rightarrow a} f(x) = \infty$. Similar considerations apply to finding a maximum.

The nature of a stationary point c can be determined by testing the second derivative $f''(c)$. If $f''(c) > 0$, then c at least qualifies as a local minimum. Similarly, if $f''(c) < 0$, then c at least qualifies as a local maximum. The indeterminate case $f''(c) = 0$ is consistent with c being a local minimum, maximum, or saddle point. For example, $f(x) = x^4$ attains its minimum at 0 while $f(x) = x^3$ has a saddle point there. In both cases, $f''(0) = 0$. Higher-order derivatives or other qualitative features of $f(x)$ must be invoked to discriminate among these possibilities. If $f''(x) \geq 0$ for all x , then $f(x)$ is said to be convex. Any stationary point of a convex function is a minimum. If $f''(x) > 0$ for all x , then $f(x)$ is strictly convex, and there is at most one stationary point. Whenever it exists, the stationary point furnishes the global minimum. A concave function satisfies $f''(x) \leq 0$ for all x . Concavity bears the same relation to maxima as convexity does to minima.

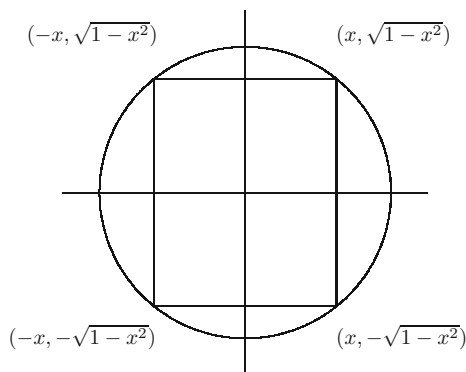
FIGURE 1.2. A *rectangle* inscribed in a *circle***Example 1.2.3 (Kepler)** *Largest Rectangle Inscribed in a Circle*

Figure 1.2 depicts a rectangle inscribed in a circle of radius 1 centered at the origin. If we suppose the vertical sides of the rectangle cross the horizontal axis at the point $(-x, 0)$ and $(x, 0)$, then Pythagoras's theorem gives the coordinates of the corners as noted in the figure. Here x is restricted to the interval $[0, 1]$. From these coordinates, it follows that the rectangle has area

$$f(x) = 4x\sqrt{1-x^2}.$$

Because $f(0) = f(1) = 0$, the maximum of $f(x)$ occurs somewhere in the open interval $(0, 1)$. Straightforward differentiation shows that

$$f'(x) = 4\sqrt{1-x^2} - \frac{4x^2}{\sqrt{1-x^2}}.$$

Setting $f'(x)$ equal to 0 and solving for x gives the critical point $x = 1/\sqrt{2}$ and the critical value $f(1/\sqrt{2}) = 2$. Since there is only one critical point on $(0, 1)$, it must be the maximum point. The largest inscribed rectangle is a square as expected. ■

Example 1.2.4 *Snell's Law*

Snell's law refers to an optical experiment involving two different media, say air and water. The less dense the medium, the faster light travels. Since light takes the path of least time, it bends at an interface such as that indicated by the horizontal axis in Fig. 1.3. Here we ask for the point $(x, 0)$ on the interface intersecting the light path. If we assume the speed of light above the interface is s_1 and below the interface is s_2 , then the total travel time is given by

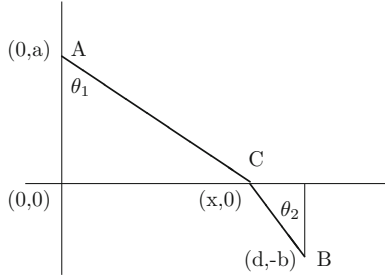


FIGURE 1.3. Diagram for Snell's law

$$f(x) = \frac{\sqrt{a^2 + x^2}}{s_1} + \frac{\sqrt{b^2 + (d-x)^2}}{s_2}. \quad (1.1)$$

The derivative of $f(x)$ is

$$f'(x) = \frac{x}{s_1\sqrt{a^2 + x^2}} - \frac{d-x}{s_2\sqrt{b^2 + (d-x)^2}}.$$

The minimum exists because $\lim_{|x| \rightarrow \infty} f(x) = \infty$. Although finding a stationary point is difficult, it is clear from the monotonicity of the functions $x/(s_1\sqrt{a^2 + x^2})$ and $(d-x)/(s_2\sqrt{b^2 + (d-x)^2})$ that it is unique. In trigonometric terms, Snell's law can be expressed as

$$\frac{\sin \theta_1}{s_1} = \frac{\sin \theta_2}{s_2}$$

using the angles at the minimum point as noted in Fig. 1.3. ■

Example 1.2.5 *The Functions* $f_n(x) = x^n e^x$

The functions $f_n(x) = x^n e^x$ for $n \geq 1$ exhibit interesting behavior. Figure 1.4 plots $f_n(x)$ for n between 1 and 3. It is clear that $\lim_{x \rightarrow -\infty} f_n(x) = 0$ and $\lim_{x \rightarrow \infty} f_n(x) = \infty$. These limits do not rule out the possibility of local maxima and minima. To find these we need

$$\begin{aligned} f'_n(x) &= (x^n + nx^{n-1})e^x \\ f''_n(x) &= [x^n + 2nx^{n-1} + n(n-1)x^{n-2}]e^x. \end{aligned}$$

Setting $f'_n(x) = 0$ produces the critical point $x = -n$, and when $n > 1$, the critical point $x = 0$. A brief calculation shows that $f''_n(-n) = (-n)^{n-1}e^{-n}$. Thus, $-n$ is a local minimum for n odd and a local maximum for n even. At 0 we have $f''_2(0) = 2$ and $f''_n(0) = 0$ for $n > 2$. Thus, the second derivative test fails for $n > 2$. However, it is clear from the variation of the sign of $f_n(x)$ to the right and left of 0 that 0 is a minimum of $f_n(x)$ for n even and

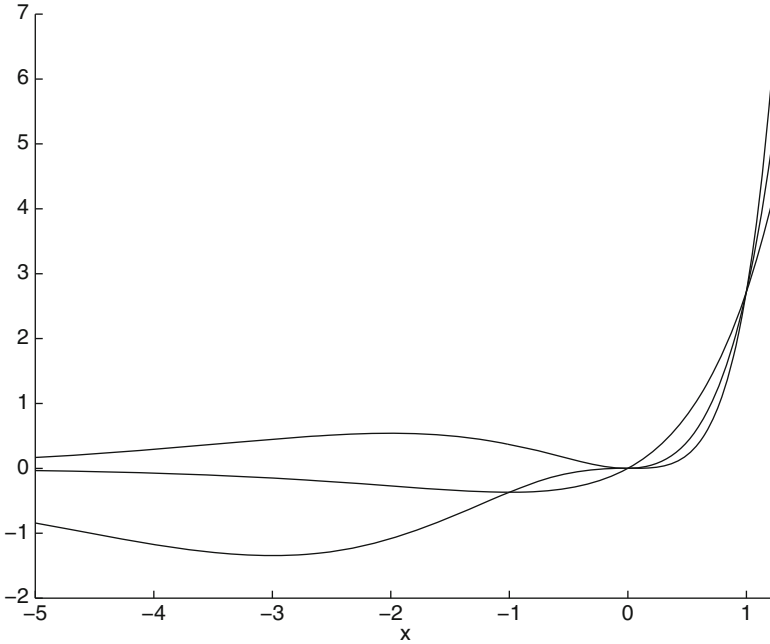


FIGURE 1.4. Plots of xe^x , x^2e^x , and x^3e^x

a saddle point of $f_n(x)$ for $n > 1$ and odd. One strength of modern graphing programs such as MATLAB is that they quickly suggest such conjectures. ■

Example 1.2.6 Fenchel Conjugate of $f_p(x) = |x|^p/p$ for $p > 1$

The Fenchel conjugate $f^*(y)$ of a convex function $f(x)$ is defined by

$$f^*(y) = \sup_x [yx - f(x)]. \tag{1.2}$$

Remarkably, $f^*(y)$ is also convex. As a particular case of this result, we consider the Fenchel conjugate of $f_p(x)$. It turns out that $f_p^*(y) = f_q(y)$, where

$$\frac{1}{p} + \frac{1}{q} = 1.$$

Here neither p nor q need be integers. According to the second derivative test, the function $f_p(x) = |x|^p/p$ is convex on the real line whenever $p > 1$. The possible failure of $f_p''(x)$ to exist at $x = 0$ does not invalidate this conclusion. To calculate $f_p^*(y)$, we observe that $f_p'(x) = |x|^{p-1} \operatorname{sgn}(x)$. This clearly implies that $x = |y|^{1/(p-1)} \operatorname{sgn}(y)$ maximizes the concave function $g(x) = yx - f_p(x)$. At the maximum point

$$\begin{aligned}
 f_p^*(y) &= yx - \frac{|x|^p}{p} \\
 &= |y|^{1+1/(p-1)} - \frac{|y|^{p/(p-1)}}{p} \\
 &= \frac{|y|^q}{q},
 \end{aligned}$$

proving our claim.

Inserting the calculated value of $f_p^*(y)$ in the definition (1.2) leads to Young's inequality

$$xy \leq \frac{|x|^p}{p} + \frac{|y|^q}{q}. \quad (1.3)$$

The double-dual identity $f_p^{**}(x) = f_p(x)$ is a special case of a general result proved later in Proposition 14.3.2. Historically, the Fenchel conjugate was introduced by Legendre for smooth functions and later generalized by Fenchel to arbitrary functions. ■

1.3 Multivariate Optimization

Although multivariate optimization is more subtle, it typically parallels univariate optimization [125, 212, 247]. The most fundamental differences arise because of constraints. In unconstrained optimization, the right definitions and notation ease the generalization. Before discussing these issues of calculus, we look at two classical inequalities that can be established by purely algebraic techniques.

Example 1.3.1 Cauchy-Schwarz Inequality

Suppose \mathbf{x} and \mathbf{y} are any two points in \mathbb{R}^n . The Cauchy-Schwarz inequality says

$$\left| \sum_{i=1}^n x_i y_i \right| \leq \left(\sum_{i=1}^n x_i^2 \right)^{1/2} \left(\sum_{i=1}^n y_i^2 \right)^{1/2}.$$

If we define the inner product

$$\mathbf{x}^* \mathbf{y} = \sum_{i=1}^n x_i y_i$$

using the transpose operator $*$ and the Euclidean norm

$$\|\mathbf{x}\| = \left(\sum_{i=1}^n x_i^2 \right)^{1/2},$$

then the inequality can be restated as $|\mathbf{x}^* \mathbf{y}| \leq \|\mathbf{x}\| \cdot \|\mathbf{y}\|$. Equality occurs in the Cauchy-Schwarz inequality if and only if \mathbf{y} is a multiple of \mathbf{x} or vice versa.

In proving the inequality, we can immediately eliminate the case $\mathbf{x} = \mathbf{0}$ where all components of \mathbf{x} are 0. Given that $\mathbf{x} \neq \mathbf{0}$, we introduce a scalar λ and consider the quadratic

$$\begin{aligned} 0 &\leq \|\lambda \mathbf{x} + \mathbf{y}\|^2 \\ &= \|\mathbf{x}\|^2 \lambda^2 + 2\mathbf{x}^* \mathbf{y} \lambda + \|\mathbf{y}\|^2 \\ &= \frac{1}{a} (a\lambda + b)^2 + c - \frac{b^2}{a} \end{aligned}$$

with $a = \|\mathbf{x}\|^2$, $b = \mathbf{x}^* \mathbf{y}$, and $c = \|\mathbf{y}\|^2$. In order for this quadratic to be nonnegative for all λ , it is necessary and sufficient that $c - b^2/a \geq 0$, which is just an abbreviation for the Cauchy-Schwarz inequality. For the quadratic to attain the value 0, the condition $c - b^2/a = 0$ must hold. When the quadratic vanishes, $\mathbf{y} = -\lambda \mathbf{x}$. ■

Example 1.3.2 Arithmetic-Geometric Mean Inequality

One generalization of the simple arithmetic-geometric mean inequality of Example 1.2.2 takes the form

$$\sqrt[n]{x_1 \cdots x_n} \leq \frac{x_1 + \cdots + x_n}{n}, \quad (1.4)$$

where x_1, \dots, x_n are any n nonnegative numbers. For a purely algebraic proof of this fact, we first note that it is obvious if any $x_i = 0$. If all $x_i > 0$, then divide both sides of the inequality by $\sqrt[n]{x_1 \cdots x_n}$. This replaces x_i by $y_i = x_i / \sqrt[n]{x_1 \cdots x_n}$ and leads to the equality $\sqrt[n]{y_1 \cdots y_n} = 1$. It now suffices to prove that $y_1 + \cdots + y_n \geq n$, which is trivially valid when $n = 1$. For $n > 1$ we argue by induction. Clearly the assumption $\sqrt[n]{y_1 \cdots y_n} = 1$ implies that there are two numbers, say y_1 and y_2 , with $y_1 \geq 1$ and $y_2 \leq 1$. If this is true, then $(y_1 - 1)(y_2 - 1) \leq 0$, or equivalently $y_1 y_2 + 1 \leq y_1 + y_2$. Invoking the induction hypothesis, we now reason that

$$\begin{aligned} y_1 + \cdots + y_n &\geq 1 + y_1 y_2 + y_3 + \cdots + y_n \\ &\geq 1 + (n - 1). \end{aligned}$$

As a prelude to discussing further examples, it is helpful to briefly summarize the theory to be developed later and often taken for granted in multidimensional calculus courses. The standard vocabulary and symbolism adopted here stress the minor adjustments necessary in going from one dimension to multiple dimensions. ■

For a real-valued function $f(\mathbf{x})$ defined on \mathbb{R}^n , the differential $df(\mathbf{x})$ is the generalization of the derivative $f'(x)$. For our purposes, $df(\mathbf{x})$ is the row vector of partial derivatives; its transpose is the gradient vector $\nabla f(\mathbf{x})$. The symmetric matrix of second partial derivatives constitutes the second differential $d^2f(\mathbf{x})$ or Hessian matrix. A stationary point \mathbf{x} satisfies $\nabla f(\mathbf{x}) = \mathbf{0}$. Fermat's principle says that all local maxima and minima on the interior of the domain of $f(\mathbf{x})$ are stationary points.

If $d^2f(\mathbf{x})$ is positive definite at a stationary point \mathbf{y} , then \mathbf{y} furnishes a local minimum. If $d^2f(\mathbf{y})$ is negative definite, then \mathbf{y} furnishes a local maximum. The function $f(\mathbf{x})$ is said to be convex if $d^2f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} ; it is strictly convex if $d^2f(\mathbf{x})$ is positive definite for all \mathbf{x} . Every stationary point of a convex function represents a global minimum. At most one stationary point exists per strictly convex function. Similar considerations apply to concave functions and global maxima, provided we substitute "negative" for "positive" throughout these definitions. These facts are rigorously proved in Chaps. 4, 5, and 6.

Example 1.3.3 *Least Squares Estimation*

Statisticians often estimate parameters by the method of least squares. To review the situation, consider n independent experiments with outcomes y_1, \dots, y_n . We wish to predict y_i from p covariates (predictors) x_{i1}, \dots, x_{ip} known in advance. For instance, y_i might be the height of the i th child in a classroom of n children. Relevant predictors might be the heights x_{i1} and x_{i2} of i 's mother and father and the sex of i coded as $x_{i3} = 1$ for a girl and $x_{i4} = 1$ for a boy. Here we take $p = 4$ and force $x_{i3}x_{i4} = 0$ so that only one sex is possible. If we use a linear predictor $\sum_{j=1}^p x_{ij}\theta_j$ of y_i , it is natural to estimate the regression coefficients θ_j by minimizing the sum of squares

$$f(\boldsymbol{\theta}) = \sum_{i=1}^n \left(y_i - \sum_{j=1}^p x_{ij}\theta_j \right)^2.$$

Differentiating $f(\boldsymbol{\theta})$ with respect to θ_j and setting the result equal to 0 produce

$$\sum_{i=1}^n x_{ij}y_i = \sum_{i=1}^n \sum_{k=1}^p x_{ij}x_{ik}\theta_k.$$

If we let \mathbf{y} denote the column vector with entries y_i and \mathbf{X} denote the matrix with entry x_{ij} in row i and column j , then these p normal equations can be written in vector form as

$$\mathbf{X}^* \mathbf{y} = \mathbf{X}^* \mathbf{X} \boldsymbol{\theta}$$

and solved as

$$\hat{\boldsymbol{\theta}} = (\mathbf{X}^* \mathbf{X})^{-1} \mathbf{X}^* \mathbf{y}.$$

In order for the indicated matrix inverse $(\mathbf{X}^* \mathbf{X})^{-1}$ to exist, $n \geq p$ should hold and the matrix \mathbf{X} must be of full rank. See Problem 22.

To check that our proposed solution $\hat{\boldsymbol{\theta}}$ represents the global minimum, we calculate the Hessian matrix $d^2 f(\boldsymbol{\theta})$. Its entries

$$\frac{\partial^2}{\partial \theta_j \partial \theta_k} f(\boldsymbol{\theta}) = 2 \sum_{i=1}^n x_{ij} x_{ik}$$

permit us to identify $d^2 f(\boldsymbol{\theta})$ with the matrix $2\mathbf{X}^* \mathbf{X}$. Owing to the full rank assumption, the symmetric matrix $\mathbf{X}^* \mathbf{X}$ is positive definite. Hence, $f(\boldsymbol{\theta})$ is strictly convex, and $\hat{\boldsymbol{\theta}}$ is the global minimum. ■

1.4 Constrained Optimization

The subject of Lagrange multipliers has a strong geometric flavor. It deals with tangent vectors and directions of steepest ascent and descent. The classical theory, which is all we consider here, is limited to equality constraints. Inequality constraints were not introduced until later in the game.

The gradient direction $\nabla f(\mathbf{x}) = df(\mathbf{x})^*$ is the direction of steepest ascent of $f(\mathbf{x})$ near the point \mathbf{x} . We can motivate this fact by considering the linear approximation

$$f(\mathbf{x} + t\mathbf{u}) = f(\mathbf{x}) + tdf(\mathbf{x})\mathbf{u} + o(t)$$

for a unit vector \mathbf{u} and a scalar t . The error term $o(t)$ becomes negligible compared to t as t decreases to 0. The inner product $df(\mathbf{x})\mathbf{u}$ in this approximation is greatest for the unit vector $\mathbf{u} = \nabla f(\mathbf{x}) / \|\nabla f(\mathbf{x})\|$. Thus, $\nabla f(\mathbf{x})$ points locally in the direction of steepest ascent of $f(\mathbf{x})$. Similarly, $-\nabla f(\mathbf{x})$ points locally in the direction of steepest descent.

Now consider minimizing or maximizing $f(\mathbf{x})$ subject to the equality constraints $g_i(\mathbf{x}) = 0$ for $i = 1, \dots, m$. A tangent direction \mathbf{w} at the point \mathbf{x} on the constraint surface satisfies $dg_i(\mathbf{x})\mathbf{w} = 0$ for all i . Of course, if the constraint surface is curved, we must interpret the tangent directions as specifying directions of infinitesimal movement. From the perpendicularity relation $dg_i(\mathbf{x})\mathbf{w} = 0$, it follows that the set of tangent directions is the orthogonal complement $S^\perp(\mathbf{x})$ of the vector subspace $S(\mathbf{x})$ spanned by the $\nabla g_i(\mathbf{x})$. To avoid degeneracies, the vectors $\nabla g_i(\mathbf{x})$ must be linearly independent. Figure 1.5 depicts level curves $g(\mathbf{x}) = c$ and gradients $\nabla g(\mathbf{x})$ for the function $\sin(x) \cos(y)$ over the square $[0, \pi] \times [-\frac{\pi}{2}, \frac{\pi}{2}]$. Tangent vectors are parallel to the level curves (contours) and perpendicular to the gradients (arrows).

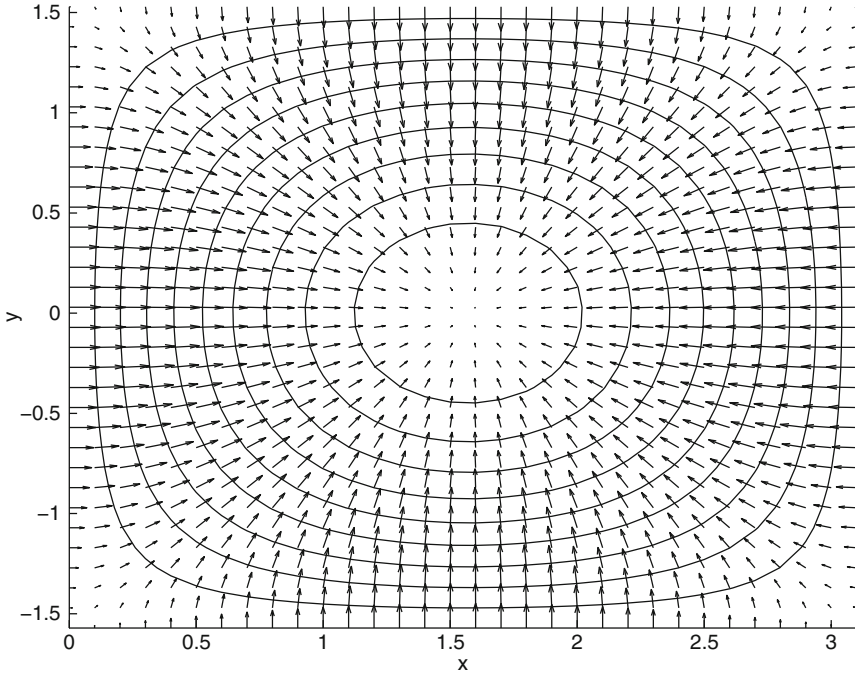


FIGURE 1.5. Level curves and steepest ascent directions for $\sin(x)\cos(y)$

At an optimal (or extremal) point \mathbf{y} , we have $df(\mathbf{y})\mathbf{w} = 0$ for every tangent direction $\mathbf{w} \in S^\perp(\mathbf{y})$; otherwise, we could move infinitesimally away from \mathbf{y} in the tangent directions \mathbf{w} and $-\mathbf{w}$ and both increase and decrease $f(\mathbf{x})$. In other words, $\nabla f(\mathbf{y})$ is a member of the double orthogonal complement $S^{\perp\perp}(\mathbf{y}) = S(\mathbf{y})$. This enables us to write

$$\nabla f(\mathbf{y}) = -\sum_{i=1}^m \lambda_i \nabla g_i(\mathbf{y})$$

for properly chosen constants $\lambda_1, \dots, \lambda_m$. Alternatively, the Lagrangian function

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\omega}) = f(\mathbf{x}) + \sum_{i=1}^m \omega_i g_i(\mathbf{x})$$

has a stationary point at $(\mathbf{y}, \boldsymbol{\lambda})$. In this regard, note that

$$\frac{\partial}{\partial \omega_i} \mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}) = 0$$

owing to the constraint $g_i(\mathbf{y}) = 0$. The essence of the Lagrange multiplier rule consists in finding the stationary points of the Lagrangian. Although

our intuitive arguments need logical tightening in many places, they offer the basic geometric insights.

Example 1.4.1 *Projection onto a Hyperplane*

A hyperplane in \mathbb{R}^n is the set of points $H = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{z}^* \mathbf{x} = c\}$ for some vector $\mathbf{z} \in \mathbb{R}^n$ and scalar c . There is no loss in generality in assuming that \mathbf{z} is a unit vector. If we seek the closest point on H to a point \mathbf{y} , then we must minimize $\|\mathbf{y} - \mathbf{x}\|^2$ subject to $\mathbf{x} \in H$. We accordingly form the Lagrangian

$$\mathcal{L}(\mathbf{x}, \lambda) = \|\mathbf{y} - \mathbf{x}\|^2 + \lambda(\mathbf{z}^* \mathbf{x} - c).$$

Setting the partial derivative with respect to x_i equal to 0 gives

$$-2(y_i - x_i) + \lambda z_i = 0.$$

This equality entails $\mathbf{x} = \mathbf{y} - \frac{1}{2}\lambda\mathbf{z}$ in vector notation. It follows that

$$c = \mathbf{z}^* \mathbf{x} = \mathbf{z}^* \mathbf{y} - \frac{1}{2}\lambda\|\mathbf{z}\|^2.$$

In view of the assumption $\|\mathbf{z}\| = 1$, we find that

$$\lambda = -2(c - \mathbf{z}^* \mathbf{y})$$

and consequently that

$$\mathbf{x} = \mathbf{y} + (c - \mathbf{z}^* \mathbf{y})\mathbf{z}.$$

If $\mathbf{y} \in H$ to begin with, then $\mathbf{x} = \mathbf{y}$. ■

Example 1.4.2 *Estimation of Multinomial Proportions*

As another statistical example, consider a multinomial experiment with m trials and observed successes m_1, \dots, m_n over n categories. The maximum likelihood estimate of the probability p_i of category i is $\hat{p}_i = m_i/m$, where $m = m_1 + \dots + m_n$. To demonstrate this fact, let

$$L(\mathbf{p}) = \binom{m}{m_1, \dots, m_n} \prod_{i=1}^n p_i^{m_i}$$

denote the likelihood. If $m_i = 0$ for some i , then we interpret $p_i^{m_i}$ as 1 even when $p_i = 0$. This convention makes it clear that we can increase $L(\mathbf{p})$ by replacing p_i by 0 and p_j by $p_j/(1 - p_i)$ for $j \neq i$. Thus, for purposes of maximum likelihood estimation, we can assume that all $m_i > 0$. Given this assumption, $L(\mathbf{p})$ tends to 0 when any p_i tends to 0. It follows that we can further restrict our attention to the interior region where all $p_i > 0$ and maximize the loglikelihood $\ln L(\mathbf{p})$ subject to the equality constraint

$\sum_{i=1}^n p_i = 1$. To find the maximum of $\ln L(\mathbf{p})$, we look for a stationary point of the Lagrangian

$$\mathcal{L}(\mathbf{p}, \lambda) = \ln \binom{m}{m_1, \dots, m_n} + \sum_{i=1}^n m_i \ln p_i + \lambda \left(\sum_{i=1}^n p_i - 1 \right).$$

Setting the partial derivative of $\mathcal{L}(\mathbf{p}, \lambda)$ with respect to p_i equal to 0 gives the equation

$$-\frac{m_i}{p_i} = \lambda.$$

These n equations are satisfied subject to the constraint by taking $\lambda = -m$ and $\hat{p}_i = m_i/m$. Thus, the necessary condition for a maximum holds at $\hat{\mathbf{p}}$. One can show that $\hat{\mathbf{p}}$ furnishes the global maximum by exploiting the strict concavity of $L(\mathbf{p})$. Although we will omit the details of this argument, it is fair to point out that strict concavity follows from

$$\frac{\partial^2}{\partial p_i \partial p_j} \ln L(\mathbf{p}) = \begin{cases} -\frac{m_i}{p_i^2} & i = j \\ 0 & i \neq j. \end{cases}$$

In statistical applications, the negative second differential $-d^2 \ln L(\mathbf{p})$ is called the observed information matrix. ■

Example 1.4.3 Eigenvalues of a Symmetric Matrix

Let $\mathbf{M} = (m_{ij})$ be an $n \times n$ symmetric matrix. Recall that \mathbf{M} has n real eigenvalues and n corresponding orthogonal eigenvectors. To find the minimum or maximum eigenvalue of \mathbf{M} , consider optimizing the quadratic form $\mathbf{x}^* \mathbf{M} \mathbf{x}$ subject to the constraint $\|\mathbf{x}\|^2 = 1$. To handle this nonlinear constraint, we introduce the Lagrangian

$$\mathcal{L}(\mathbf{x}, \lambda) = \mathbf{x}^* \mathbf{M} \mathbf{x} + \lambda(1 - \|\mathbf{x}\|^2).$$

Setting the partial derivative of $\mathcal{L}(\mathbf{x}, \lambda)$ with respect to x_i equal to 0 yields

$$2 \sum_{j=1}^n m_{ij} x_j - 2\lambda x_i = 0.$$

In matrix notation, this reduces to $\mathbf{M} \mathbf{x} = \lambda \mathbf{x}$. It follows that

$$\mathbf{x}^* \mathbf{M} \mathbf{x} = \lambda \mathbf{x}^* \mathbf{x} = \lambda.$$

Thus, the stationary points of the Lagrangian are eigenvectors of \mathbf{M} . The Lagrange multipliers are the corresponding stationary values or eigenvalues. The maximum and minimum eigenvalues occur among these stationary values. This result does not directly invoke the spectral decomposition theorem.

To prove that every symmetric matrix possesses a spectral decomposition, one can argue by induction. The scalar case is trivial, so suppose that the claim is true for every $(n - 1) \times (n - 1)$ symmetric matrix. If \mathbf{M} is an $n \times n$ symmetric matrix, then as just noted, \mathbf{M} has a maximum eigenvalue λ and corresponding unit eigenvector \mathbf{x} . Consider the deflated matrix $\mathbf{N} = \mathbf{M} - \lambda\mathbf{x}\mathbf{x}^*$ and the subspace $S = \{c\mathbf{x} : c \in \mathbb{R}\}$. It is clear that \mathbf{N} maps S into $\mathbf{0}$. Furthermore, \mathbf{N} also maps the perpendicular complement S^\perp of S into itself. Indeed, if $\mathbf{y} \in S^\perp$, then

$$\mathbf{x}^* \mathbf{N} \mathbf{y} = \mathbf{x}^* (\mathbf{M} \mathbf{y} - \lambda \mathbf{x} \mathbf{x}^* \mathbf{y}) = \lambda \mathbf{x}^* \mathbf{y} - \lambda \mathbf{x}^* \mathbf{y} = 0.$$

By the induction hypothesis, the symmetric linear transformation induced by \mathbf{N} on S^\perp possesses a spectral decomposition. The $n - 1$ eigenvectors of this decomposition, which are automatically perpendicular to \mathbf{x} , also serve as eigenvectors of \mathbf{M} . ■

Example 1.4.4 A Minimization Problem with Two Constraints

In three dimensions the plane $x_1 + x_2 + x_3 = 1$ intersects the cylinder $x_1^2 + x_2^2 = 1$ in an ellipse. To find the closest point on the ellipse to the origin, we construct the Lagrangian

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = x_1^2 + x_2^2 + x_3^2 + \lambda_1(x_1 + x_2 + x_3 - 1) + \lambda_2(x_1^2 + x_2^2 - 1).$$

The stationary points of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ satisfy the system of equations

$$\begin{aligned} 2x_1 + \lambda_1 + 2\lambda_2 x_1 &= 0 \\ 2x_2 + \lambda_1 + 2\lambda_2 x_2 &= 0 \\ 2x_3 + \lambda_1 &= 0 \end{aligned} \tag{1.5}$$

in addition to the constraints. Since $\lambda_1 = -2x_3$, the first two equations of the system (1.5) can be recast as

$$\begin{aligned} (1 + \lambda_2)x_1 &= x_3 \\ (1 + \lambda_2)x_2 &= x_3. \end{aligned}$$

If $\lambda_2 = -1$, then $x_3 = 0$. In this case it is geometrically obvious that $(1, 0, 0)$ and $(0, 1, 0)$ are the only two points that satisfy the constraints $x_1 + x_2 = 1$ and $x_1^2 + x_2^2 = 1$. On the other hand, if $\lambda_2 \neq -1$, then $x_1 = x_2$. In this case, the constraints $2x_1 + x_3 = 1$ and $2x_1^2 = 1$ dictate the solutions $(\frac{\sqrt{2}}{2}, \frac{\sqrt{2}}{2}, 1 - \sqrt{2})$ and $(-\frac{\sqrt{2}}{2}, -\frac{\sqrt{2}}{2}, 1 + \sqrt{2})$. Of the four candidate points, it is easy to check that $(1, 0, 0)$ and $(0, 1, 0)$ are closest to the origin. ■

Example 1.4.5 A Population Genetics Problem

The multiplier rule is sometimes hard to apply, and ad hoc methods can lead to better results. In the setting of the multinomial distribution, consider the problem of maximizing the sum

$$f(\mathbf{p}) = \sum_{i < j} (2p_i p_j)^2 \frac{1}{2}$$

subject to the constraints $\sum_{i=1}^n p_i = 1$ and all $p_i \geq 0$. This problem has a genetics interpretation involving a locus with n codominant alleles labeled $1, \dots, n$. (See Sect. 8.4 for some genetics terminology.) At the locus, most genotype combinations of a mother, father, and child make it possible to infer which allele the mother contributes to the child and which allele the father contributes to the child. It turns out that the only ambiguous case occurs when all three family members share the same heterozygous genotype i/j , where $i \neq j$. The probability of this configuration is $(2p_i p_j)^2 \frac{1}{2}$ if p_i and p_j are the population frequencies (proportions) of alleles i and j . Here $2p_i p_j$ is the frequency of an i/j mother or an i/j father and $\frac{1}{2}$ is the probability that one of them transmits an i allele and the other transmits a j allele. Thus, $f(\mathbf{p})$ represents the probability that the trio's genotypes do not permit inference of the child's maternal and paternal alleles.

The case $n = 2$ is particularly simple because the function $f(\mathbf{p})$ then reduces to $2(p_1 p_2)^2$. In view of Example 1.2.2, the maximum of $\frac{1}{8}$ is attained when $p_1 = p_2 = \frac{1}{2}$. This suggests that the maximum for general n occurs when all $p_i = \frac{1}{n}$. Because there are $\binom{n}{2}$ heterozygous genotypes,

$$f\left(\frac{1}{n}\mathbf{1}\right) = \binom{n}{2} \left(\frac{2}{n^2}\right)^2 \frac{1}{2} = \frac{n-1}{n^3},$$

which is strictly less than $\frac{1}{8}$ for $n \geq 3$. Our first guess is wrong, and we now conjecture that the maximum occurs on a boundary where all but two of the $p_i = 0$. If we permute the components of a maximum point, then symmetry dictates that the result will also be a maximum point. We therefore order the parameters so that $0 < p_1 \leq p_2 \leq \dots \leq p_n$, avoiding for the moment the lower-dimensional case where $p_1 = 0$.

We now argue that we can increase $f(\mathbf{p})$ by increasing p_2 by $q \in [0, p_1]$ at the expense of decreasing p_1 by q . Consider the function

$$g(q) = 2(p_1 - q)^2 \sum_{i=3}^n p_i^2 + 2(p_2 + q)^2 \sum_{i=3}^n p_i^2 + 2(p_1 - q)^2 (p_2 + q)^2$$

which equals the original objective function except for an additive constant independent of q . For $n \geq 3$, straightforward differentiation gives

$$\begin{aligned} g'(q) &= -4(p_1 - q) \sum_{i=3}^n p_i^2 + 4(p_2 + q) \sum_{i=3}^n p_i^2 \\ &\quad - 4(p_1 - q)(p_2 + q)^2 + 4(p_1 - q)^2 (p_2 + q) \end{aligned}$$

$$\begin{aligned}
&= 4(p_2 - p_1 + 2q) \left[\sum_{i=3}^n p_i^2 - (p_1 - q)(p_2 + q) \right] \\
&\geq 4(p_2 - p_1 + 2q) \left[\sum_{i=3}^n p_i^2 - (p_2 - q)(p_2 + q) \right] \\
&= 4(p_2 - p_1 + 2q) \left[\sum_{i=4}^n p_i^2 + p_3^2 - p_2^2 + q^2 \right] \\
&\geq 0
\end{aligned}$$

for $q \in [0, p_1]$. Thus, we should reduce p_1 to 0 and increase p_2 to $p_2 + p_1$. Furthermore, we should keep discarding the lowest positive p_j until all but two of the p_i equal 0. Finally, we set the remaining two p_i equal to $\frac{1}{2}$. This verifies our second conjecture. ■

Example 1.4.6 *Polygon of Greatest Area Inscribed in a Circle*

As a generalization of Example 1.2.3, consider a polygon inscribed in a circle. For a given number of vertices, the polygon with the greatest area is regular. Here is a proof using elementary plane geometry [213]. Let i , j , and k be three successive vertices as depicted in Fig. 1.6. If we fix the positions of i and k , then we can ask for the optimal placement of vertex j . The only part of the polygon in play is the triangle ijk . The area of the triangle equals half its base times its height. The base distance from i to k is fixed, and the height is clearly a maximum when j is moved to the position j' on the perpendicular bisector of the base. When $j = j'$, the sides ij and jk have equal length. Repeating this argument for all successive vertex triples demonstrates that all sides must have equal length and that the polygon is regular. This is a constrained optimization problem because all vertices are forced to lie on the given circle. Presumably one could reach the conclusion that the polygon is regular by invoking Lagrange multipliers, but the necessary machinery would obscure the underlying simplicity of the problem. ■

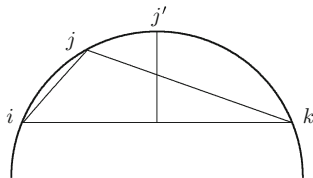


FIGURE 1.6. The *triangle* defined by three adjacent vertices

1.5 Problems

1. Given a point C in the interior of an acute angle, find the points A and B on the sides of the angle such that the perimeter of the triangle ABC is as short as possible. (Hint: Reflect C perpendicularly across the two sides of the angle to points C_1 and C_2 , respectively. Let A and B be the intersections of the line segment connecting C_1 and C_2 with the two sides.)
2. Find the point in a triangle that minimizes the sum of the squared distances from the vertices. Show that this point is the intersection of the medians of the triangle.
3. Given an angle in the plane and a point in its interior, find the line that passes through the point and cuts off from the angle a triangle of minimal area. This triangle is determined by the vertex of the angle and the two points where the constructed line intersects the sides of the angle.
4. Find the minima of the functions

$$\begin{aligned} f(x) &= x \ln x \\ g(x) &= x - \ln x \\ h(x) &= x + \frac{1}{x} \end{aligned}$$

on $(0, \infty)$. Demonstrate rigorously that your solutions are indeed the minima.

5. For $t > 0$ prove that $e^x > x^t$ for all $x > 0$ if and only if $t < e$ [69].
6. Consider the function

$$f(x) = 2(x+2)(x \ln x - x + 1) - 3(x-1)^2$$

defined on the interval $(0, \infty)$. Show that $f(x) \geq f(1) = 0$ for all x . (Hint: Expand $f(x)$ in a third-order Taylor series around $x = 1$.)

7. Demonstrate that Euler's function $f(x) = x^2 - 1/\ln x$ possesses no local or global minima on either domain $(0, 1)$ or $(1, \infty)$.
8. Prove the harmonic-geometric mean inequality

$$\frac{1}{\frac{1}{n} \left(\frac{1}{x_1} + \cdots + \frac{1}{x_n} \right)} \leq \sqrt[n]{x_1 \cdots x_n}$$

for n positive numbers x_1, \dots, x_n .

9. Prove the arithmetic-quadratic mean inequality

$$\frac{1}{n} \sum_{i=1}^n x_i \leq \left(\frac{1}{n} \sum_{i=1}^n x_i^2 \right)^{1/2}$$

for any nonnegative numbers x_1, \dots, x_n .

10. Heron's classical formula for the area of a triangle with sides of length a , b , and c is $\sqrt{s(s-a)(s-b)(s-c)}$, where $s = (a+b+c)/2$ is the semiperimeter. Show that the triangle of fixed perimeter with greatest area is equilateral.

11. Consider the sequences

$$e_n = \left(1 + \frac{1}{n}\right)^n, \quad f_n = \left(1 - \frac{1}{n}\right)^n.$$

It is well known that $\lim_{n \rightarrow \infty} e_n = e$ and $\lim_{n \rightarrow \infty} f_n = e^{-1}$. Use the arithmetic-geometric mean inequality to prove that $e_n \leq e_{n+1}$ and $f_n \leq f_{n+1}$. In addition, prove that $e_n \leq f_n^{-1}$, that e_n and f_n^{-1} have finite limits, and that these limits are equal [34]. (Hint: Write e_n as the product of 1 and n copies of $1 + \frac{1}{n}$.)

12. Demonstrate that the function

$$f(\mathbf{x}) = 4x_1 + \frac{x_1}{x_2^2} + \frac{4x_2}{x_1}$$

on \mathbb{R}^2 has the minimum value 8 for x_1 and x_2 positive. At what point \mathbf{x} is the minimum attained? (Hint: Write

$$f(\mathbf{x}) = 4 \left(\frac{4x_1 + \frac{x_1}{x_2^2} + \frac{2x_2}{x_1} + \frac{2x_2}{x_1}}{4} \right)$$

and apply the arithmetic-geometric mean inequality [34]. Attacking this problem by calculus is harder.)

13. Let $H_n = 1 + \frac{1}{2} + \dots + \frac{1}{n}$. Verify the inequality $n \sqrt[n]{n+1} \leq n + H_n$ for any positive integer n (Putnam Competition, 1975).
14. Consider an n -gon circumscribing the unit circle in \mathbb{R}^2 . Demonstrate that the n -gon has minimum area if and only if all of its n sides have equal length. (Hint: Let θ_m be the circular angle between the two points of tangency of sides m and $m+1$ [69]. Show that the area of the quadrilateral defined by the center of the circle, the two points of tangency, and the intersection of the two sides is given by $\tan \frac{\theta_m}{2}$.)

15. Find the minimum of the function

$$g(\mathbf{x}) = 2x_1^2 + x_2^2 + \frac{1}{2x_1^2 + x_2^2}$$

on \mathbb{R}^2 . (Hint: Consider $f(x) = x + 1/x$ on $(0, \infty)$.)

16. In forensic applications of genetics, the sum

$$s(\mathbf{p}) = 1 - 2\left(\sum_{i=1}^n p_i^2\right)^2 + \sum_{i=1}^n p_i^4$$

occurs [165]. Here the p_i are nonnegative and sum to 1. Prove rigorously that $s(\mathbf{p})$ attains its maximum $s_{\max} = 1 - \frac{2}{n^2} + \frac{1}{n^3}$ when all $p_i = \frac{1}{n}$. (Hint: To prove the claim about s_{\max} , note that without loss of generality one can assume $p_1 \leq p_2 \leq \dots \leq p_n$. If $p_i < p_{i+1}$, then $s(\mathbf{p})$ can be increased by replacing p_i and p_{i+1} by $p_i + x$ and $p_{i+1} - x$ for x positive and sufficiently small.)

17. Suppose that a and b are real numbers satisfying $0 < a < b$. Prove that the origin locally minimizes $f(\mathbf{x}) = (x_2 - ax_1^2)(x_2 - bx_1^2)$ along every line $x_1 = ht$ and $x_2 = kt$ through the origin. Also show that $f(t, ct^2) < 0$ for $a < c < b$ and $t \neq 0$. The origin therefore affords a local minimum along each line through the origin but not a local minimum in the wider sense. If $c < a$ or $c > b$, then $f(t, ct^2) > 0$ for $t \neq 0$, and the paradox disappears.

18. Demonstrate that the function $x_1^2 + x_2^2(1 - x_1)^3$ has a unique stationary point in \mathbb{R}^2 , which is a local minimum but not a global minimum. Can this occur for a continuously differentiable function with domain \mathbb{R}^2 ?

19. Find all of the stationary points of the function

$$f(\mathbf{x}) = x_1^2 x_2 e^{-x_1^2 - x_2^2}$$

in \mathbb{R}^2 . Classify each point as either a local minimum, a local maximum, or a saddle point.

20. Rosenbrock's function $f(\mathbf{x}) = 100(x_1^2 - x_2)^2 + (x_1 - 1)^2$ achieves its global minimum at $\mathbf{x} = \mathbf{1}$. Prove that $\nabla f(\mathbf{1}) = \mathbf{0}$ and $d^2 f(\mathbf{1})$ is positive definite.

21. Consider the polynomial

$$p(\mathbf{x}) = \left[x_1^2 + (1 - x_2)^2\right] \left[x_2^2 + (1 - x_1)^2\right]$$

in two variables. Show that $p(\mathbf{x})$ is symmetric in the sense that $p(x_1, x_2) = p(x_2, x_1)$, that $\lim_{\|\mathbf{x}\| \rightarrow \infty} p(\mathbf{x}) = \infty$, and that $p(\mathbf{x})$ does not attain its minimum along the diagonal $x_1 = x_2$.

22. Suppose that the $m \times n$ matrix \mathbf{X} has full rank and that $m \geq n$. Show that the $n \times n$ matrix $\mathbf{X}^* \mathbf{X}$ is invertible and positive definite.
23. Consider two sets of positive numbers x_1, \dots, x_n and $\alpha_1, \dots, \alpha_n$ such that $\sum_{i=1}^n \alpha_i = 1$. Prove the generalized arithmetic-geometric mean inequality

$$\prod_{i=1}^n x_i^{\alpha_i} \leq \sum_{i=1}^n \alpha_i x_i$$

by minimizing $\sum_{i=1}^n \alpha_i x_i$ subject to the constraint $\prod_{i=1}^n x_i^{\alpha_i} = c$.

24. Suppose the components of a vector $\mathbf{x} \in \mathbb{R}^n$ are positive and have product $\prod_{k=1}^n x_k = 1$. Prove that

$$\prod_{k=1}^n (1 + x_k) \geq 2^n.$$

25. Find the rectangular box in \mathbb{R}^3 of greatest volume having a fixed surface area.
26. Let $S(\mathbf{0}, r) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| = r\}$ be the sphere of radius r centered at the origin. For $\mathbf{y} \in \mathbb{R}^n$, find the point of $S(\mathbf{0}, r)$ closest to \mathbf{y} .
27. Find the parallelepiped of maximum volume that can be inscribed in the ellipsoid

$$\frac{x_1^2}{a^2} + \frac{x_2^2}{b^2} + \frac{x_3^2}{c^2} = 1.$$

Assume that the parallelepiped is centered at the origin and has edges parallel to the coordinate axes.

28. A twice continuously differentiable function $f(\mathbf{x})$ on \mathbb{R}^2 satisfies

$$\frac{\partial^2}{\partial x_1^2} f(\mathbf{x}) + \frac{\partial^2}{\partial x_2^2} f(\mathbf{x}) > 0$$

for all \mathbf{x} . Prove that $f(\mathbf{x})$ has no local maxima [69]. An example of such a function is $f(\mathbf{x}) = \|\mathbf{x}\|^2 = x_1^2 + x_2^2$.

29. Use the Cauchy-Schwarz inequality to verify the inequalities

$$\sum_{m=0}^n a_m x^m \leq \frac{1}{\sqrt{1-x^2}} \left(\sum_{m=0}^n a_m^2 \right)^{\frac{1}{2}} \quad 0 \leq x < 1$$

$$\sum_{m=1}^n \frac{a_m}{m} \leq \sqrt{\frac{\pi^2}{6}} \left(\sum_{m=1}^n a_m^2 \right)^{\frac{1}{2}}$$

$$\sum_{m=1}^n \frac{a_m}{\sqrt{m+n}} \leq \sqrt{\ln 2} \left(\sum_{m=1}^n a_m^2 \right)^{\frac{1}{2}}$$

$$\sum_{m=0}^n \binom{n}{m} a_m \leq \binom{2n}{n}^{\frac{1}{2}} \left(\sum_{m=1}^n a_m^2 \right)^{\frac{1}{2}}.$$

The upper bound n can be finite or infinite in the first two cases [243].

30. Verify Lagrange's identity

$$\left(\sum_{i=1}^n x_i y_i \right)^2 = \sum_{i=1}^n x_i^2 \sum_{j=1}^n y_j^2 - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n (x_i y_j - x_j y_i)^2.$$

How does this lead to a proof of the Cauchy-Schwarz inequality and its stated condition for equality?

31. Demonstrate the bound

$$\left| \sum_{m=1}^n a_m \right|^2 + \left| \sum_{m=1}^n (-1)^m a_m \right|^2 \leq (n+2) \sum_{m=1}^n a_m^2.$$

This is better than the obvious bound $2n \sum_{m=1}^n a_m^2$ given by the Cauchy-Schwarz inequality [243]. (Hint: Let e_n and o_n be the sum of the a_m with m even and odd, respectively.)

32. Consider the function $f(x) = \sum_{m=1}^n p_m \cos(\alpha_m x)$ for a discrete probability distribution p_1, \dots, p_n . Given that $g(x) = \cos \alpha x$ satisfies the identity $g(x)^2 = \frac{1}{2}[1 + g(2x)]$, show that $f(x)$ satisfies the inequality $f(x)^2 \leq \frac{1}{2}[1 + f(2x)]$. This is the Harker-Kasper inequality from X-ray crystallography [243].
33. For positive numbers b_1, \dots, b_n and h_1, \dots, h_n , show that

$$\min_{1 \leq m \leq n} \frac{h_m}{b_m} \leq \frac{h_1 + \dots + h_n}{b_1 + \dots + b_n} \leq \max_{1 \leq m \leq n} \frac{h_m}{b_m}.$$

This inequality of Cauchy has the baseball interpretation that a batting average of a team is never worse than that of its worst player and never better than that of its best player [243]. (Hint: Consider the case $n = 2$ and use induction.)

2

The Seven C's of Analysis

2.1 Introduction

The current chapter explains key concepts of mathematical analysis summarized by the six adjectives convergent, complete, closed, compact, continuous, and connected. Chapter 6 will add to these six c's the seventh c, convex. At first blush these concepts seem remote from practical problems of optimization. However, painful experience and exotic counterexamples have taught mathematicians to pay attention to details. Fortunately, we can benefit from the struggles of earlier generations and bypass many of the intellectual traps.

2.2 Vector and Matrix Norms

In multidimensional calculus, vector and matrix norms quantify notions of topology and convergence [48, 105, 117, 207]. Norms are also helpful in estimating rates of convergence of iterative methods for solving linear and nonlinear equations and optimizing functions. Functional analysis, which deals with infinite-dimensional vector spaces, uses norms on functions.

We have already met the Euclidean vector norm $\|\mathbf{x}\|$ on \mathbb{R}^n . For most purposes, this norm suffices. It shares with other norms the four properties:

- (a) $\|\mathbf{x}\| \geq 0$,
- (b) $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$,

(c) $\|c\mathbf{x}\| = |c| \cdot \|\mathbf{x}\|$ for every real number c ,

(d) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.

Property (d) is known as the triangle inequality. To prove it for the Euclidean norm, we note that the Cauchy-Schwarz inequality implies

$$\begin{aligned}\|\mathbf{x} + \mathbf{y}\|^2 &= \|\mathbf{x}\|^2 + 2\mathbf{x}^*\mathbf{y} + \|\mathbf{y}\|^2 \\ &\leq \|\mathbf{x}\|^2 + 2\|\mathbf{x}\|\|\mathbf{y}\| + \|\mathbf{y}\|^2 \\ &= (\|\mathbf{x}\| + \|\mathbf{y}\|)^2.\end{aligned}$$

One immediate consequence of the triangle inequality is the further inequality

$$\left| \|\mathbf{x}\| - \|\mathbf{y}\| \right| \leq \|\mathbf{x} - \mathbf{y}\|.$$

Two other simple but helpful norms are the ℓ_1 and ℓ_∞ norms

$$\begin{aligned}\|\mathbf{x}\|_1 &= \sum_{i=1}^n |x_i| \\ \|\mathbf{x}\|_\infty &= \max_{1 \leq i \leq n} |x_i|.\end{aligned}$$

Some of the properties of these norms are explored in the problems. In the mathematical literature, the three norms are often referred to as the ℓ_2 , ℓ_1 , and ℓ_∞ norms.

An $m \times n$ matrix $\mathbf{A} = (a_{ij})$ can be viewed as a vector in \mathbb{R}^{mn} . Accordingly, we define its Frobenius norm

$$\|\mathbf{A}\|_F = \left(\sum_{i=1}^m \sum_{j=1}^n a_{ij}^2 \right)^{1/2} = \sqrt{\text{tr}(\mathbf{A}\mathbf{A}^*)} = \sqrt{\text{tr}(\mathbf{A}^*\mathbf{A})},$$

where $\text{tr}(\cdot)$ is the matrix trace function. Our reasons for writing $\|\mathbf{A}\|_F$ rather than $\|\mathbf{A}\|$ will soon be apparent. In the meanwhile, the Frobenius matrix norm satisfies the additional condition

(e) $\|\mathbf{A}\mathbf{B}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{B}\|$

for any two compatible matrices $\mathbf{A} = (a_{ij})$ and $\mathbf{B} = (b_{ij})$. Property (e) is verified by invoking the Cauchy-Schwarz inequality in

$$\begin{aligned}\|\mathbf{A}\mathbf{B}\|_F^2 &= \sum_{i,j} \left| \sum_k a_{ik}b_{kj} \right|^2 \\ &\leq \sum_{i,j} \left(\sum_k a_{ik}^2 \right) \left(\sum_l b_{lj}^2 \right) \\ &= \left(\sum_{i,k} a_{ik}^2 \right) \left(\sum_{l,j} b_{lj}^2 \right) \\ &= \|\mathbf{A}\|_F^2 \|\mathbf{B}\|_F^2.\end{aligned}\tag{2.1}$$

The Frobenius norm does not satisfy the natural condition

$$(f) \|\mathbf{I}\| = 1$$

for an identity matrix \mathbf{I} . Indeed, an easy calculation shows that $\|\mathbf{I}\|_F = \sqrt{n}$ when \mathbf{I} is $n \times n$.

To meet all of the conditions (a) through (f), we need to turn to induced matrix norms. Let $\|\cdot\|$ denote both the Euclidean norm on \mathbb{R}^m and the Euclidean norm on \mathbb{R}^n . The induced Euclidean norm on $m \times n$ matrices is defined by

$$\|\mathbf{A}\| = \sup_{\mathbf{x} \neq \mathbf{0}} \frac{\|\mathbf{A}\mathbf{x}\|}{\|\mathbf{x}\|} = \sup_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{x}\|. \quad (2.2)$$

For reasons explained in Proposition 2.2.1, the induced norm (2.2) is called the spectral norm. The question of whether the indicated supremum exists definition (2.2) is settled by the inequalities

$$\|\mathbf{A}\mathbf{x}\| \leq \sum_{i=1}^n |x_i| \cdot \|\mathbf{A}\mathbf{e}_i\| \leq \left(\sum_{i=1}^n \|\mathbf{A}\mathbf{e}_i\| \right) \|\mathbf{x}\|,$$

where $\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i$ and \mathbf{e}_i is the unit vector whose entries are all 0 except for $e_{ii} = 1$. More exotic induced matrix norms can be concocted by substituting non-Euclidean norms in the numerator and denominator of definition (2.2). For square matrices, the two norms ordinarily coincide. All of the defining properties of a matrix norm are trivial to check for an induced matrix norm. For instance, property (e) follows from

$$\begin{aligned} \|\mathbf{A}\mathbf{B}\| &= \sup_{\|\mathbf{x}\|=1} \|\mathbf{A}\mathbf{B}\mathbf{x}\| \\ &\leq \|\mathbf{A}\| \sup_{\|\mathbf{x}\|=1} \|\mathbf{B}\mathbf{x}\| \\ &= \|\mathbf{A}\| \cdot \|\mathbf{B}\|. \end{aligned}$$

Definition (2.2) also clearly entails the equality $\|\mathbf{I}\| = 1$ when $m = n$.

The next proposition determines the value of the Euclidean norm $\|\mathbf{A}\|$. In the proposition, $\rho(\mathbf{M})$ denotes the absolute value of the dominant eigenvalue of the square matrix \mathbf{M} . This quantity is called the spectral radius of \mathbf{M} .

Proposition 2.2.1 *If $\mathbf{A} = (a_{ij})$ is an $m \times n$ matrix, then*

$$\|\mathbf{A}\| = \sqrt{\rho(\mathbf{A}^* \mathbf{A})} = \sqrt{\rho(\mathbf{A} \mathbf{A}^*)} = \|\mathbf{A}^*\|.$$

When \mathbf{A} is symmetric, $\|\mathbf{A}\|$ reduces to $\rho(\mathbf{A})$. The norms $\|\mathbf{A}\|$ and $\|\mathbf{A}\|_F$ satisfy

$$\|\mathbf{A}\| \leq \|\mathbf{A}\|_F \leq \sqrt{n} \|\mathbf{A}\|. \quad (2.3)$$

Finally, when \mathbf{A} is a row or column vector, the Euclidean matrix and vector norms of \mathbf{A} coincide.

Proof: Choose an orthonormal basis of eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ for the symmetric matrix $\mathbf{A}^* \mathbf{A}$ with corresponding eigenvalues arranged so that $0 \leq \lambda_1 \leq \dots \leq \lambda_n$. If $\mathbf{x} = \sum_{i=1}^n c_i \mathbf{u}_i$ is a unit vector, then $\sum_{i=1}^n c_i^2 = 1$, and

$$\begin{aligned} \|\mathbf{A}\|^2 &= \sup_{\|\mathbf{x}\|=1} \mathbf{x}^* \mathbf{A}^* \mathbf{A} \mathbf{x} \\ &= \sup_{\|\mathbf{x}\|=1} \sum_{i=1}^n \lambda_i c_i^2 \\ &\leq \lambda_n. \end{aligned}$$

Equality is achieved when $c_n = \pm 1$ and all other $c_i = 0$. If \mathbf{A} is symmetric with eigenvalues μ_i arranged so that $|\mu_1| \leq \dots \leq |\mu_n|$, then the \mathbf{u}_i can be chosen to be the corresponding eigenvectors. In this case, clearly $\lambda_i = \mu_i^2$.

To prove that $\rho(\mathbf{A}^* \mathbf{A}) = \rho(\mathbf{A} \mathbf{A}^*)$, choose an eigenvalue $\lambda \neq 0$ of $\mathbf{A}^* \mathbf{A}$ with corresponding eigenvector \mathbf{v} . Multiplying the equation $\mathbf{A}^* \mathbf{A} \mathbf{v} = \lambda \mathbf{v}$ on the left by \mathbf{A} produces $(\mathbf{A} \mathbf{A}^*) \mathbf{A} \mathbf{v} = \lambda \mathbf{A} \mathbf{v}$. Because $\mathbf{A}^* \mathbf{A} \mathbf{v} = \lambda \mathbf{v}$, the vector $\mathbf{A} \mathbf{v} \neq \mathbf{0}$. Thus, λ is an eigenvalue of $\mathbf{A} \mathbf{A}^*$ with eigenvector $\mathbf{A} \mathbf{v}$. Likewise, any eigenvalue $\omega \neq 0$ of $\mathbf{A} \mathbf{A}^*$ is also an eigenvalue of $\mathbf{A}^* \mathbf{A}$.

To verify the left bound of the pair of bounds (2.3), apply inequality (2.1) with $\mathbf{B} = \mathbf{x}$ in the definition of $\|\mathbf{A}\|$. The right bound follows from

$$\sum_{i=1}^m a_{ij}^2 = \|\mathbf{A} \mathbf{e}_j\|^2 \leq \|\mathbf{A}\|^2$$

by summing on j . Finally, suppose that \mathbf{A} is a column vector. The two bounds (2.3) with $n = 1$ show that $\|\mathbf{A}\| = \|\mathbf{A}\|_F$. If \mathbf{A} is a row vector, the same reasoning applied to \mathbf{A}^* gives $\|\mathbf{A}\| = \|\mathbf{A}^*\| = \|\mathbf{A}^*\|_F = \|\mathbf{A}\|_F$. ■

2.3 Convergence and Completeness

A sequence $\mathbf{x}_m \in \mathbb{R}^n$ converges to \mathbf{x} , written $\lim_{m \rightarrow \infty} \mathbf{x}_m = \mathbf{x}$, provided $\lim_{m \rightarrow \infty} \|\mathbf{x}_m - \mathbf{x}\| = 0$ in the standard Euclidean norm. For convergence of \mathbf{x}_m to \mathbf{x} to occur, it is necessary and sufficient that each component sequence x_{mi} converge to x_i . Convergence of a sequence of matrices is defined similarly using either the Frobenius norm $\|\mathbf{A}\|_F$ or the induced matrix norm $\|\mathbf{A}\|$. The pair of bounds (2.3) shows that the two norms are equivalent in testing convergence.

Convergent sequences of vectors or matrices enjoy many useful properties. Some of these are mentioned in the next proposition.

Proposition 2.3.1 *In the following list, once a limit is assumed to exist for an item, it is assumed to exist for all subsequent items. With this proviso, we have:*

(a) *If $\lim_{m \rightarrow \infty} \mathbf{x}_m = \mathbf{x}$, then $\lim_{m \rightarrow \infty} \|\mathbf{x}_m\| = \|\mathbf{x}\|$.*

(b) If $\lim_{m \rightarrow \infty} \mathbf{y}_m = \mathbf{y}$, then

$$\lim_{m \rightarrow \infty} \mathbf{x}_m^* \mathbf{y}_m = \mathbf{x}^* \mathbf{y}.$$

(c) If a and b are real scalars, then

$$\lim_{m \rightarrow \infty} [a\mathbf{x}_m + b\mathbf{y}_m] = a\mathbf{x} + b\mathbf{y}.$$

(d) If $\lim_{m \rightarrow \infty} \mathbf{M}_m = \mathbf{M}$ for a sequence of matrices compatible with \mathbf{x} , then

$$\lim_{m \rightarrow \infty} \mathbf{M}_m \mathbf{x}_m = \mathbf{M} \mathbf{x}.$$

(e) If \mathbf{M} is square and invertible, then \mathbf{M}_m^{-1} exists for large m and

$$\lim_{m \rightarrow \infty} \mathbf{M}_m^{-1} = \mathbf{M}^{-1}.$$

(f) Finally, if $\lim_{m \rightarrow \infty} \mathbf{N}_m = \mathbf{N}$ for a sequence of matrices compatible with \mathbf{M} , then

$$\lim_{m \rightarrow \infty} \mathbf{M}_m \mathbf{N}_m = \mathbf{M} \mathbf{N}.$$

Proof: As a sample proof, part (d) follows from the inequalities

$$\begin{aligned} \|\mathbf{M}_m \mathbf{x}_m - \mathbf{M} \mathbf{x}\| &\leq \|\mathbf{M}_m \mathbf{x}_m - \mathbf{M}_m \mathbf{x}\| + \|\mathbf{M}_m \mathbf{x} - \mathbf{M} \mathbf{x}\| \\ &\leq \|\mathbf{M}_m\| \cdot \|\mathbf{x}_m - \mathbf{x}\| + \|\mathbf{M}_m - \mathbf{M}\| \cdot \|\mathbf{x}\| \\ \|\mathbf{M}_m\| &\leq \|\mathbf{M}_m - \mathbf{M}\| + \|\mathbf{M}\|. \end{aligned}$$

Part (e) will be proved after Example 2.3.3. ■

In some situations, we know that the members of a sequence become progressively closer together. A Cauchy sequence \mathbf{x}_m exhibits a strong form of this phenomenon; namely, for every $\epsilon > 0$, there is an m such that $\|\mathbf{x}_p - \mathbf{x}_q\| \leq \epsilon$ for all $p, q \geq m$. The real line \mathbb{R} is complete in the sense that every Cauchy sequence possesses a limit. The rational numbers are incomplete by contrast because a sequence of rationals can converge to an irrational. The completeness of \mathbb{R} carries over to \mathbb{R}^n . Indeed, if \mathbf{x}_m is a Cauchy sequence, then under the Euclidean norm we have

$$|x_{pi} - x_{qi}| \leq \|\mathbf{x}_p - \mathbf{x}_q\|.$$

This shows that each component sequence is Cauchy and consequently possesses a limit x_i . The vector \mathbf{x} with components x_i then furnishes a limit for the vector sequence \mathbf{x}_m .

Example 2.3.1 *Existence of Suprema and Infima*

The completeness of the real line is equivalent to the existence of least upper bounds or suprema. Consider a nonempty set $S \subset \mathbb{R}$ that is bounded above. If the set is finite, then its least upper bound is just its largest element. If the set is infinite, we choose a and b such that the interval $[a, b]$ contains an element of S and b is an upper bound of S . We can generate $\sup S$ by a bisection strategy. Bisect $[a, b]$ into the two subintervals $[a, (a+b)/2]$ and $[(a+b)/2, b]$. Let $[a_1, b_1]$ denote the left subinterval if $(a+b)/2$ provides an upper bound. Otherwise, let $[a_1, b_1]$ denote the right subinterval. In either case, $[a_1, b_1]$ contains an element of S . Now bisect $[a_1, b_1]$ and generate a subinterval $[a_2, b_2]$ by the same criterion. If we continue bisecting and choosing a left or right subinterval ad infinitum, then we generate two Cauchy sequences a_i and b_i with common limit c . By the definition of the sequence b_i , c furnishes an upper bound of S . By the definition of the sequence a_i , no bound of S is smaller than c . Establishing the existence of the greatest lower bound $\inf S$ for S bounded below proceeds similarly. If S is unbounded above, then $\sup S = \infty$, and if it is unbounded below, then $\inf S = -\infty$. ■

Example 2.3.2 *Limit Superior and Limit Inferior*

For a real sequence x_n , we define the limit superior and limit inferior by

$$\begin{aligned}\limsup_{n \rightarrow \infty} x_n &= \inf_m \sup_{n \geq m} x_n = \lim_{m \rightarrow \infty} \sup_{n \geq m} x_n \\ \liminf_{n \rightarrow \infty} x_n &= \sup_m \inf_{n \geq m} x_n = \lim_{m \rightarrow \infty} \inf_{n \geq m} x_n.\end{aligned}$$

If $\sup_n x_n = \infty$, then $\limsup_{n \rightarrow \infty} x_n = \infty$, and if $\lim_{n \rightarrow \infty} x_n = -\infty$, then $\limsup_{n \rightarrow \infty} x_n = -\infty$. From these definitions, one can also deduce that

$$\limsup_{n \rightarrow \infty} -x_n = -\liminf_{n \rightarrow \infty} x_n \quad (2.4)$$

and that

$$\liminf_{n \rightarrow \infty} x_n \leq \limsup_{n \rightarrow \infty} x_n. \quad (2.5)$$

The sequence x_n has a limit if and only if equality prevails in inequality (2.5). In this situation, the common value of the limit superior and inferior furnishes the limit of x_n . ■

Example 2.3.3 *Series Expansion for a Matrix Inverse*

If a square matrix M has norm $\|M\| < 1$, then we can write

$$(\mathbf{I} - M)^{-1} = \sum_{i=0}^{\infty} M^i.$$

To verify this claim, we first prove that the partial sums $\mathbf{S}_j = \sum_{i=0}^j \mathbf{M}^i$ form a Cauchy sequence. This fact is a consequence of the inequalities

$$\begin{aligned} \|\mathbf{S}_k - \mathbf{S}_j\| &= \left\| \sum_{i=j+1}^k \mathbf{M}^i \right\| \\ &\leq \sum_{i=j+1}^k \|\mathbf{M}^i\| \\ &\leq \sum_{i=j+1}^k \|\mathbf{M}\|^i \end{aligned}$$

for $k \geq j$ and the assumption $\|\mathbf{M}\| < 1$. If we let \mathbf{S} represent the limit of the \mathbf{S}_j , then part (f) of Proposition 2.3.1 implies that $(\mathbf{I} - \mathbf{M})\mathbf{S}_j$ converges to $(\mathbf{I} - \mathbf{M})\mathbf{S}$. But $(\mathbf{I} - \mathbf{M})\mathbf{S}_j = \mathbf{I} - \mathbf{M}^{j+1}$ also converges to \mathbf{I} . Hence, $(\mathbf{I} - \mathbf{M})\mathbf{S} = \mathbf{I}$, and this verifies the claim $\mathbf{S} = (\mathbf{I} - \mathbf{M})^{-1}$. ■

With this result under our belts, we now demonstrate part (e) of Proposition 2.3.1. Because $\|\mathbf{M}^{-1}(\mathbf{M} - \mathbf{M}_m)\| \leq \|\mathbf{M}^{-1}\| \cdot \|\mathbf{M} - \mathbf{M}_m\|$, the matrix inverse $[\mathbf{I} - \mathbf{M}^{-1}(\mathbf{M} - \mathbf{M}_m)]^{-1}$ exists for large m . Therefore, we can write the inverse of

$$\begin{aligned} \mathbf{M}_m &= \mathbf{M} - (\mathbf{M} - \mathbf{M}_m) \\ &= \mathbf{M}[\mathbf{I} - \mathbf{M}^{-1}(\mathbf{M} - \mathbf{M}_m)] \end{aligned}$$

as

$$\mathbf{M}_m^{-1} = [\mathbf{I} - \mathbf{M}^{-1}(\mathbf{M} - \mathbf{M}_m)]^{-1} \mathbf{M}^{-1}.$$

The proof of convergence is completed by noting the bound

$$\begin{aligned} \|\mathbf{M}_m^{-1} - \mathbf{M}^{-1}\| &= \left\| \sum_{i=1}^{\infty} [\mathbf{M}^{-1}(\mathbf{M} - \mathbf{M}_m)]^i \mathbf{M}^{-1} \right\| \\ &\leq \sum_{i=1}^{\infty} \|\mathbf{M}^{-1}\|^i \|\mathbf{M} - \mathbf{M}_m\|^i \|\mathbf{M}^{-1}\| \\ &= \frac{\|\mathbf{M}^{-1}\|^2 \|\mathbf{M} - \mathbf{M}_m\|}{1 - \|\mathbf{M}^{-1}\| \cdot \|\mathbf{M} - \mathbf{M}_m\|}, \end{aligned}$$

applying in the process the matrix analog of part (a) of Proposition 2.3.1.

Example 2.3.4 Matrix Exponential Function

The exponential of a square matrix \mathbf{M} is given by the series expansion

$$e^{\mathbf{M}} = \sum_{i=0}^{\infty} \frac{1}{i!} \mathbf{M}^i.$$

To prove the convergence of the series, it again suffices to show that the partial sums $\mathbf{S}_j = \sum_{i=0}^j \frac{1}{i!} \mathbf{M}^i$ form a Cauchy sequence. The bound

$$\|\mathbf{S}_k - \mathbf{S}_j\| = \left\| \sum_{i=j+1}^k \frac{1}{i!} \mathbf{M}^i \right\| \leq \sum_{i=j+1}^k \frac{1}{i!} \|\mathbf{M}\|^i$$

for $k \geq j$ is just what we need.

The matrix exponential function has many interesting properties. For example, the function $N(t) = e^{t\mathbf{M}}$ solves the differential equation

$$N'(t) = \mathbf{M}N(t)$$

subject to the initial condition $N(0) = \mathbf{I}$. Here t is a real parameter, and we differentiate the matrix $N(t)$ entry by entry. In Example 4.2.2 of Chap. 4, we will prove that $N(t) = e^{t\mathbf{M}}$ is the one and only solution. The law of exponents $e^{\mathbf{A}+\mathbf{B}} = e^{\mathbf{A}}e^{\mathbf{B}}$ for commuting matrices \mathbf{A} and \mathbf{B} is another interesting property of the matrix exponential function. One way of proving the law of exponents is to observe that $e^{t(\mathbf{A}+\mathbf{B})}$ and $e^{t\mathbf{A}}e^{t\mathbf{B}}$ both solve the differential equation

$$N'(t) = (\mathbf{A} + \mathbf{B})N(t)$$

subject to the initial condition $N(0) = \mathbf{I}$. Since the solution to such an initial value problem is unique, the two solutions must coincide at $t = 1$. ■

2.4 The Topology of \mathbb{R}^n

Mathematics involves a constant interplay between the abstract and the concrete. We now consider some qualitative features of sets in \mathbb{R}^n that generalize to more abstract spaces. For instance, there is the matter of boundedness. A set $S \subset \mathbb{R}^n$ is said to be bounded if it is contained in some ball $B(\mathbf{0}, r) = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < r\}$ of radius r centered at the origin $\mathbf{0}$. As we shall see in our discussion of compactness, boundedness takes on added importance when it is combined with the notion of closedness. A closed set is closed under the formation of limits. Thus, $S \subset \mathbb{R}^n$ is closed if for every convergent sequence \mathbf{x}_m taken from S , we have $\lim_{m \rightarrow \infty} \mathbf{x}_m \in S$ as well.

It takes time and effort to appreciate the ramifications of these ideas. A few of the most pertinent ones for closedness are noted in the next proposition.

Proposition 2.4.1 *The collection of closed sets satisfy the following:*

- (a) *The whole space \mathbb{R}^n is closed.*
- (b) *The empty set \emptyset is closed.*

- (c) The intersection $S = \bigcap_{\alpha} S_{\alpha}$ of an arbitrary number of closed sets S_{α} is closed.
- (d) The union $S = \bigcup_{\alpha} S_{\alpha}$ of a finite number of closed sets S_{α} is closed.

Proof: All of these are easy. For part (d), observe that for any convergent sequence \mathbf{x}_m taken from S , one of the sets S_{α} must contain an infinite subsequence \mathbf{x}_{m_k} . The limit of this subsequence exists and falls in S_{α} . ■

Some examples of closed sets are closed intervals $(-\infty, a]$, $[a, b]$, and $[b, \infty)$; closed balls $\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| \leq r\}$; spheres $\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| = r\}$; hyperplanes $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{z}^* \mathbf{x} = c\}$; and closed halfspaces $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{z}^* \mathbf{x} \leq c\}$. A closed set S of \mathbb{R}^n is complete in the sense that all Cauchy sequences from S possess limits in S .

Example 2.4.1 Finitely Generated Convex Cones

A set C is a convex cone provided $\alpha \mathbf{u} + \beta \mathbf{v}$ is in C whenever the vectors \mathbf{u} and \mathbf{v} are in C and the scalars α and β are nonnegative. A finitely generated convex cone can be written as

$$C = \left\{ \sum_{i=1}^m \alpha_i \mathbf{v}_i : \alpha_i \geq 0, i = 1, \dots, m \right\}.$$

Demonstrating that C is a closed set is rather subtle. Consider a sequence $\mathbf{u}_j = \sum_{i=1}^m \alpha_{ji} \mathbf{v}_i$ in C converging to a point \mathbf{u} . If the vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ are linearly independent, then the coefficients α_{ji} are the unique coordinates of \mathbf{u}_j in the finite-dimensional subspace spanned by the \mathbf{v}_i . To recover the α_{ji} , we introduce the matrix \mathbf{V} with columns $\mathbf{v}_1, \dots, \mathbf{v}_m$ and rewrite the original equation as $\mathbf{u}_j = \mathbf{V} \boldsymbol{\alpha}_j$. Multiplying this equation by first \mathbf{V}^* and then by $(\mathbf{V}^* \mathbf{V})^{-1}$ on the left gives $\boldsymbol{\alpha}_j = (\mathbf{V}^* \mathbf{V})^{-1} \mathbf{V}^* \mathbf{u}_j$. This representations allows us to conclude that $\boldsymbol{\alpha}_j$ possesses a limit $\boldsymbol{\alpha}$ with nonnegative entries. Therefore, the limit $\mathbf{u} = \mathbf{V} \boldsymbol{\alpha}$ lies in C .

If we relax the assumption that the vectors are linearly independent, we must resort to an inductive argument to prove that C is closed. The case $m = 1$ is true because a single vector \mathbf{v}_1 is linearly independent. Assume that the claim holds for $m - 1$ vectors. If the vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ are linearly independent, then we are done. If the vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ are linearly dependent, then there exist scalars β_1, \dots, β_m , not all 0, such that $\sum_{i=1}^m \beta_i \mathbf{v}_i = \mathbf{0}$. Without loss of generality, we can assume that $\beta_i < 0$ for at least one index i . We can express any point $\mathbf{u} \in C$ as

$$\mathbf{u} = \sum_{i=1}^m \alpha_i \mathbf{v}_i = \sum_{i=1}^m (\alpha_i + t \beta_i) \mathbf{v}_i$$

for an arbitrary scalar t . If we increase t gradually from 0, then there is a first value at which $\alpha_j + t \beta_j = 0$ for some index j . This shows that C can

be decomposed as the union

$$C = \bigcup_{j=1}^m \left\{ \sum_{i \neq j} \gamma_i \mathbf{v}_i : \gamma_i \geq 0, i \neq j \right\}.$$

Each of the convex cones $\{\sum_{i \neq j} \gamma_i \mathbf{v}_i : \gamma_i \geq 0, i \neq j\}$ is closed by the induction hypothesis. Since a finite union of closed sets is closed, C itself is closed. A straightforward extension of this argument establishes the stronger claim that every point in the cone can be represented as a positive combination of a linearly independent subset of $\{\mathbf{v}_1, \dots, \mathbf{v}_m\}$. ■

The complement $S^c = \mathbb{R}^n \setminus S$ of a closed set S is called an open set. Every $\mathbf{x} \in S^c$ is surrounded by a ball $B(\mathbf{x}, r)$ completely contained in S^c . If this were not the case, then we could construct a sequence of points \mathbf{x}_m from S converging to \mathbf{x} , contradicting the closedness of S . This fact is the first of several mentioned in the next proposition.

Proposition 2.4.2 *The collection of open sets satisfy the following:*

- (a) *Every open set is a union of balls, and every union of balls is an open set.*
- (b) *The whole space \mathbb{R}^n is open.*
- (c) *The empty set \emptyset is open.*
- (d) *The union $S = \cup_{\alpha} S_{\alpha}$ of an arbitrary number of open sets S_{α} is open.*
- (e) *The intersection $S = \cap_{\alpha} S_{\alpha}$ of a finite number of open sets S_{α} is open.*

Proof: Again these are easy. Parts (d) and (e) are consequences of the set identities

$$\begin{aligned} (\cap_{\alpha} S_{\alpha})^c &= \cup_{\alpha} S_{\alpha}^c \\ (\cup_{\alpha} S_{\alpha})^c &= \cap_{\alpha} S_{\alpha}^c \end{aligned}$$

and parts (c) and (d) of Proposition 2.4.1. ■

Some examples of open sets are open intervals $(-\infty, a)$, (a, b) , and (b, ∞) ; balls $\{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| < r\}$, and open halfspaces $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{z}^* \mathbf{x} < c\}$. Any open set surrounding a point is called a neighborhood of the point. Some examples of sets that are neither closed nor open are the unbalanced intervals $(a, b]$ and $[a, b)$, the discrete set $V = \{n^{-1} : n = 1, 2, \dots\}$, and the rational numbers. If we append the limit 0 to the set V , then it becomes closed.

A boundary point \mathbf{x} of a set S is the limit of a sequence of points from S and also the limit of a different sequence of points from S^c . Closed sets contain all of their boundary points, and open sets contain none of their

boundary points. The interior of S is the largest open set contained within S . The closure of S is the smallest closed set containing S . For instance, the boundary of the ball $B(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| < r\}$ is the sphere $S(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| = r\}$. The closure of $B(\mathbf{x}, r)$ is the closed ball $C(\mathbf{x}, r) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{x}\| \leq r\}$, and the interior of $C(\mathbf{x}, r)$ is $B(\mathbf{x}, r)$.

A closed bounded set is said to be compact. Finite intervals $[a, b]$ are typical compact sets. Compact sets can be defined in several equivalent ways. The most important of these is the Bolzano-Weierstrass characterization. In preparation for this result, let us define a multidimensional interval $[\mathbf{a}, \mathbf{b}]$ in \mathbb{R}^n to be the Cartesian product

$$[\mathbf{a}, \mathbf{b}] = [a_1, b_1] \times \cdots \times [a_n, b_n]$$

of n one-dimensional intervals. We will only consider closed intervals. The diameter of $[\mathbf{a}, \mathbf{b}]$ is the greatest separation between any two of its points; this clearly reduces to the distance $\|\mathbf{a} - \mathbf{b}\|$ between its extreme corners.

Proposition 2.4.3 (Bolzano-Weierstrass) *A set $S \subset \mathbb{R}^n$ is compact if and only if every sequence \mathbf{x}_m in S has a convergent subsequence \mathbf{x}_{m_i} with limit in S .*

Proof: Suppose every sequence \mathbf{x}_m in S has a convergent subsequence \mathbf{x}_{m_i} with limit in S . If S is unbounded, then we can define a sequence \mathbf{x}_m with $\|\mathbf{x}_m\| \geq m$. Clearly, this sequence has no convergent subsequence. If S is not closed, then there is a convergent sequence \mathbf{x}_m with limit \mathbf{x} outside S . Clearly, no subsequence of \mathbf{x}_m can converge to a limit in S . Thus, the subsequence property implies compactness.

For the converse, let \mathbf{x}_m be a sequence in the compact set S . Because S is bounded, it is contained in a multidimensional interval $[\mathbf{a}, \mathbf{b}]$. If infinitely many of the \mathbf{x}_m coincide, then these can be used to construct a constant subsequence that trivially converges to a point of S . Otherwise, let T_0 denote the infinite set $\cup_{m=1}^{\infty} \{\mathbf{x}_m\}$.

The rest of the proof adapts the bisection strategy of Example 2.3.1. The first stage of the bisection divides $[\mathbf{a}, \mathbf{b}]$ into 2^n subintervals of equal volume. Each of these subintervals can be written as $[\mathbf{a}_1, \mathbf{b}_1]$, where $a_{1j} = a_j$ and $b_{1j} = (a_j + b_j)/2$ or $a_{1j} = (a_j + b_j)/2$ and $b_{1j} = b_j$. There is no harm in the fact that these subintervals overlap along their boundaries. It is only vital to observe that one of the subintervals contains an infinite subset $T_1 \subset T_0$. Let us choose such a subinterval and label it using the generic notation $[\mathbf{a}_1, \mathbf{b}_1]$. We now inductively repeat the process. At stage $i + 1$ we divide the previously chosen subinterval $[\mathbf{a}_i, \mathbf{b}_i]$ into 2^n subintervals of equal volume. Each of these subintervals can be written as $[\mathbf{a}_{i+1}, \mathbf{b}_{i+1}]$, where either $a_{i+1,j} = a_{ij}$ and $b_{i+1,j} = (a_{ij} + b_{ij})/2$ or $a_{i+1,j} = (a_{ij} + b_{ij})/2$ and $b_{i+1,j} = b_{ij}$. One of these subintervals, which we label $[\mathbf{a}_{i+1}, \mathbf{b}_{i+1}]$ for convenience, contains an infinite subset $T_{i+1} \subset T_i$.

We continue this process ad infinitum. In the process choosing \mathbf{x}_{m_i} from T_i and $m_i > m_{i-1}$. Because $T_i \subset [\mathbf{a}_i, \mathbf{b}_i]$ and the diameter of $[\mathbf{a}_i, \mathbf{b}_i]$ tends

to 0, the subsequence \mathbf{x}_{m_i} is Cauchy. By virtue of the completeness of \mathbb{R}^n , this subsequence converges to some point \mathbf{x} , which necessarily belongs to the closed set S . ■

In many instances it is natural to consider a subset S of \mathbb{R}^n as a topological space in its own right. Notions of distance and convergence carry over immediately, but we must exercise some care in defining closed and open sets. In the relative topology, a subset $T \subset S$ is closed if and only if it can be represented as the intersection $T = S \cap C$ of S with a closed set C of \mathbb{R}^n . If T is closed in S , then the obvious choice of C is the closure of T in \mathbb{R}^n . Likewise, $T \subset S$ is open in the relative topology if and only if it can be represented as the intersection $T = S \cap O$ of S with an open set O of \mathbb{R}^n . These two definitions are consistent with an open set being the relative complement of a closed set and vice versa. They are also consistent with the development of continuous functions sketched in the next section.

2.5 Continuous Functions

Continuous functions are the building blocks of mathematical analysis. Continuity is such an intuitive notion that ancient mathematicians did not even bother to define it. Proper recognition of continuity had to wait until differentiability was thoroughly explored. Our approach to continuity emphasizes convergent sequences. A function $f(\mathbf{x})$ from \mathbb{R}^m to \mathbb{R}^n is said to be continuous at \mathbf{y} if $f(\mathbf{x}_i)$ converges to $f(\mathbf{y})$ for every sequence \mathbf{x}_i that converges to \mathbf{y} . If the domain of $f(\mathbf{x})$ is a subset S of \mathbb{R}^m , then the sequences \mathbf{x}_i and the point \mathbf{y} are confined to S . Finally, $f(\mathbf{x})$ is said to be continuous if it is continuous at every point \mathbf{y} of its domain.

The definition of continuity through convergent sequences tends to be simpler to apply than the competing ϵ and δ approach of calculus. We leave it to the reader to show that the two definitions are fully equivalent. Either definition has powerful consequences. For instance, it is clear that a vector-valued function is continuous if and only if each of its component functions is continuous. Before enumerating other less obvious consequences, it is helpful to forge a few tools for recognizing and constructing continuous functions. Fortunately, the collection of continuous functions is closed under many standard algebraic operations. Here are a few examples.

Proposition 2.5.1 *Given that the vector-valued functions $f(\mathbf{x})$ and $g(\mathbf{x})$ and matrix-valued function $M(\mathbf{x})$ and $N(\mathbf{x})$ are continuous and compatible whenever necessary, the following algebraic combinations are continuous:*

- (a) The norm $\|f(\mathbf{x})\|$.
- (b) The inner product $f(\mathbf{x})^*g(\mathbf{x})$.
- (c) The linear combination $\alpha f(\mathbf{x}) + \beta g(\mathbf{x})$ for real scalars α and β .

- (d) The matrix-vector product $M(\mathbf{x})f(\mathbf{x})$.
- (e) The matrix inverse $M^{-1}(\mathbf{x})$ when $M(\mathbf{x})$ is square and invertible.
- (f) The matrix product $M(\mathbf{x})N(\mathbf{x})$.
- (g) The functional composition $f \circ g(\mathbf{x}) = f[g(\mathbf{x})]$.

Proof: Parts (a) through (f) are all immediate by-products of Proposition 2.3.1 and the definition of continuity. For part (g), suppose \mathbf{x}_i tends to \mathbf{x} . Then $f(\mathbf{x}_i)$ tends to $f(\mathbf{x})$, and so $f \circ g(\mathbf{x}_i)$ tends to $f \circ g(\mathbf{x})$. ■

Example 2.5.1 Rational Functions

Because the coordinate variables x_i of $\mathbf{x} \in \mathbb{R}^n$ are continuous, all polynomials in these variables are continuous as well. For example, the determinant of a square matrix is a continuous function of the entries of the matrix. A quotient of two polynomials (rational function) in the coordinate variables x_i of $\mathbf{x} \in \mathbb{R}^n$ is continuous where its denominator does not vanish. Finally, any linear transformation of one vector space into another is continuous. ■

Example 2.5.2 Distance to a Set

The distance $\text{dist}(\mathbf{x}, S)$ from a point $\mathbf{x} \in \mathbb{R}^n$ to a set S is defined by

$$\text{dist}(\mathbf{x}, S) = \inf_{\mathbf{z} \in S} \|\mathbf{z} - \mathbf{x}\|.$$

To prove that the function $\text{dist}(\mathbf{x}, S)$ is continuous in \mathbf{x} , take the infimum over $\mathbf{z} \in S$ of both sides of the triangle inequality

$$\|\mathbf{z} - \mathbf{x}\| \leq \|\mathbf{z} - \mathbf{y}\| + \|\mathbf{y} - \mathbf{x}\|.$$

This demonstrates that $\text{dist}(\mathbf{x}, S) \leq \text{dist}(\mathbf{y}, S) + \|\mathbf{y} - \mathbf{x}\|$. Reversing the roles of \mathbf{x} and \mathbf{y} then leads to the inequality

$$|\text{dist}(\mathbf{x}, S) - \text{dist}(\mathbf{y}, S)| \leq \|\mathbf{y} - \mathbf{x}\|,$$

establishing continuity. ■

In generalizing continuity to more abstract topological spaces, the characterizations in the next proposition are crucial.

Proposition 2.5.2 *The following conditions are equivalent for a function $f(\mathbf{x})$ from $T \subset \mathbb{R}^m$ to \mathbb{R}^n :*

- (a) $f(\mathbf{x})$ is continuous.
- (b) The inverse image $f^{-1}(S)$ of every closed set S is closed.
- (c) The inverse image $f^{-1}(S)$ of every open set S is open.

Proof: To prove that (a) implies (b), suppose \mathbf{x}_i is a sequence in $f^{-1}(S)$ tending to $\mathbf{x} \in T$. Then the conclusion $\lim_{i \rightarrow \infty} f(\mathbf{x}_i) = f(\mathbf{x})$ identifies $f(\mathbf{x})$ as an element of the closed set S and therefore \mathbf{x} as belonging to $f^{-1}(S)$. Conditions (b) and (c) are equivalent because of the relation

$$f^{-1}(S)^c = f^{-1}(S^c)$$

between inverse images and set complements. Finally, to prove that (c) entails (a), suppose that $\lim_{i \rightarrow \infty} \mathbf{x}_i = \mathbf{x}$. For any $\epsilon > 0$, the inverse image of the ball $B[f(\mathbf{x}), \epsilon]$ is open by assumption. Consequently, there exists a neighborhood $T \cap B(\mathbf{x}, \delta)$ mapped into $B[f(\mathbf{x}), \epsilon]$. In other words,

$$\|f(\mathbf{x}_i) - f(\mathbf{x})\| < \epsilon$$

whenever $\|\mathbf{x}_i - \mathbf{x}\| < \delta$, which is sufficient to validate continuity. ■

Example 2.5.3 *Continuity of $\sqrt[m]{x}$*

The root function $f(x) = \sqrt[m]{x}$ is the functional inverse of the power function $g(x) = x^m$. We have already noted that $g(x)$ is continuous. On the interval $(0, \infty)$, it is also strictly increasing and maps the open interval (a, b) onto the open interval (a^m, b^m) . Put another way, $f^{-1}[(a, b)] = (a^m, b^m)$. (Here we implicitly invoke the intermediate value property proved in Proposition 2.7.1.) Because the inverse image of a union of open intervals is a union of open intervals, application of part (c) of Proposition 2.5.2 establishes the continuity of $f(x)$. ■

Example 2.5.4 *The Set of Positive Definite Matrices*

A symmetric $n \times n$ matrix $\mathbf{M} = (m_{ij})$ can be viewed as a point in \mathbb{R}^m for $m = \binom{n}{2} + n$. To demonstrate that the subset S of positive definite matrices is open in \mathbb{R}^m , we invoke the classical criterion of Sylvester. (See Problem 29 of Chap. 5 or [136].) This test for positive definiteness uses the determinants of the principal submatrices of \mathbf{M} . The k th of these submatrices \mathbf{M}_k is the $k \times k$ upper left block of \mathbf{M} . If \mathbf{M} is positive definite, then one can show that \mathbf{M}_k is positive definite by taking a nontrivial $k \times 1$ vector \mathbf{x}_k and padding it with zeros to construct a nontrivial $n \times 1$ vector \mathbf{x} . It is then clear that $\mathbf{x}_k^* \mathbf{M}_k \mathbf{x}_k = \mathbf{x}^* \mathbf{M} \mathbf{x} > 0$. Because \mathbf{M}_k is positive definite, its determinant $\det \mathbf{M}_k > 0$. Conversely, if all of the $\det \mathbf{M}_k > 0$, then \mathbf{M} itself is positive definite.

Given this background, we write

$$S = \bigcap_{k=1}^n \{\mathbf{M} : \det \mathbf{M}_k > 0\}.$$

Because the functions $\det \mathbf{M}_k$ are continuous in the entries of \mathbf{M} , the inverse images $\{\mathbf{M} : \det \mathbf{M}_k > 0\}$ of the open set $(0, \infty)$ are open. Since a finite intersection of open sets is open, S itself is an open set. ■

As opposed to inverse images, the image of a closed (open) set under a continuous function need not be closed (open). However, continuous functions do preserve compactness.

Proposition 2.5.3 *Suppose the continuous function $f(\mathbf{x})$ maps the compact set $S \subset \mathbb{R}^m$ into \mathbb{R}^n . Then the image $f(S)$ is compact.*

Proof: The key is to apply Proposition 2.4.3. Let $f(\mathbf{x}_i)$ be a sequence in $f(S)$. Extract a convergent subsequence \mathbf{x}_{i_j} of \mathbf{x}_i with limit $\mathbf{y} \in S$. Then the continuity of $f(\mathbf{x})$ compels $f(\mathbf{x}_{i_j})$ to converge to $f(\mathbf{y})$. ■

We now come to one of the most important results in optimization theory.

Proposition 2.5.4 (Weierstrass) *Let $f(\mathbf{x})$ be a continuous real-valued function defined on a set S of \mathbb{R}^n . If the set $T = \{\mathbf{x} \in S : f(\mathbf{x}) \geq f(\mathbf{y})\}$ is compact for some $\mathbf{y} \in S$, then $f(\mathbf{x})$ attains its supremum on S . Similarly, if $T = \{\mathbf{x} \in S : f(\mathbf{x}) \leq f(\mathbf{y})\}$ is compact for some $\mathbf{y} \in S$, then $f(\mathbf{x})$ attains its infimum on S . Both conclusions apply when S itself is compact.*

Proof: Consider the question of whether the function $f(\mathbf{x})$ attains its supremum $u = \sup_{\mathbf{x} \in S} f(\mathbf{x})$. The set $f(T)$ is bounded by virtue of Proposition 2.5.3, and the supremum of $f(\mathbf{x})$ on T coincides with u . For every positive integer i choose a point $\mathbf{x}_i \in T$ such that $f(\mathbf{x}_i) \geq u - 1/i$. In view of the compactness of T , we can extract a convergent subsequence of \mathbf{x}_i with limit $\mathbf{z} \in T$. The continuity of $f(\mathbf{x})$ along this subsequence then implies that $f(\mathbf{z}) = u$. ■

Example 2.5.5 *Closest Point in a Set*

To prove that the distance $\text{dist}(\mathbf{x}, S)$ is achieved for some $\mathbf{z} \in S$, we must assume that S is closed. In finding the closest point to \mathbf{x} in S , choose any point $\mathbf{y} \in S$. The set $T = S \cap \{\mathbf{z} : \|\mathbf{z} - \mathbf{x}\| \leq \|\mathbf{y} - \mathbf{x}\|\}$ is both closed and bounded and therefore compact. Proposition 2.5.4 now informs us that the continuous function $\mathbf{z} \mapsto \|\mathbf{z} - \mathbf{x}\|$ attains its infimum on S . ■

Example 2.5.6 *Equivalence of Norms*

Every norm $\|\mathbf{x}\|_{\dagger}$ on \mathbb{R}^n is equivalent to the Euclidean norm $\|\mathbf{x}\|$ in the sense that there exist positive constants a and b such that the inequalities

$$a\|\mathbf{x}\| \leq \|\mathbf{x}\|_{\dagger} \leq b\|\mathbf{x}\| \quad (2.6)$$

hold for all \mathbf{x} . To prove the right inequality in (2.6), let $\mathbf{e}_1, \dots, \mathbf{e}_n$ denote the standard basis. Then conditions (c) and (d) defining a norm indicate that $\mathbf{x} = \sum_{i=1}^n x_i \mathbf{e}_i$ satisfies

$$\|\mathbf{x}\|_{\dagger} \leq \sum_{i=1}^n |x_i| \cdot \|\mathbf{e}_i\|_{\dagger} = \|\mathbf{x}\| \sum_{i=1}^n \|\mathbf{e}_i\|_{\dagger}.$$

This proves the upper bound with $b = \sum_{i=1}^n \|\mathbf{e}_i\|_{\dagger}$.

To establish the lower bound, we note that property (c) of a norm allows us to restrict attention to the sphere $S = \{\mathbf{x} : \|\mathbf{x}\| = 1\}$. Now the function $\mathbf{x} \mapsto \|\mathbf{x}\|_{\dagger}$ is uniformly continuous on \mathbb{R}^n because

$$|\|\mathbf{x}\|_{\dagger} - \|\mathbf{y}\|_{\dagger}| \leq \|\mathbf{x} - \mathbf{y}\|_{\dagger} \leq b\|\mathbf{x} - \mathbf{y}\|$$

follows from the upper bound just demonstrated. Since the sphere S is compact, the continuous function $\mathbf{x} \rightarrow \|\mathbf{x}\|_{\dagger}$ attains its lower bound a on S . In view of property (b) defining a norm, $a > 0$. ■

Example 2.5.7 *The Fundamental Theorem of Algebra*

Consider a polynomial $p(z) = c_n z^n + c_{n-1} z^{n-1} + \cdots + c_0$ in the complex variable z with $c_n \neq 0$. The fundamental theorem of algebra says that $p(z)$ has a root. d'Alembert suggested an interesting optimization proof of this fact [30]. We begin by observing that if we identify a complex number with an ordered pair of real numbers, then the domain of the real-valued function $|p(z)|$ is \mathbb{R}^2 . The identity

$$|p(z)| = |z|^n \left| c_n + \frac{c_{n-1}}{z} + \cdots + \frac{c_0}{z^n} \right|$$

shows that $|p(z)|$ tends to ∞ whenever $|z|$ tends to ∞ . Therefore, the set $T = \{z : |p(z)| \leq d\}$ is compact for any d , and Proposition 2.5.4 implies that $|p(z)|$ attains its minimum at some point y . Expanding $p(z)$ around y gives a polynomial

$$q(z) = p(z + y) = b_n z^n + b_{n-1} z^{n-1} + \cdots + b_0$$

with the same degree as $p(z)$. Furthermore, the minimum of $|q(z)|$ occurs at $z = 0$. Suppose $b_1 = \cdots = b_{k-1} = 0$ and $b_k \neq 0$. For some angle $\theta \in [0, 2\pi)$, the scaled complex exponential

$$u = \left| \frac{b_0}{b_k} \right|^{1/k} e^{i\theta/k}$$

is a root of the equation $b_k u^k + b_0 = 0$. The function $f(t) = |q(tu)|$ clearly satisfies $f(t) \geq |b_0|$ and

$$f(t) = |b_k t^k u^k + b_0| + o(t^k) = |b_0(1 - t^k)| + o(t^k)$$

for t small and positive. These two conditions are compatible only if $b_0 = 0$. Hence, the minimum of $|q(z)| = |p(z + y)|$ is 0. ■

Example 2.5.8 *Continuity of the Roots of a Polynomial*

As a followup to the previous example, let us prove that the roots of a polynomial depend continuously on its coefficients [261]. One has to exercise

caution in stating this result. First, we limit ourselves to monic polynomials $p(z) = z^n + c_{n-1}z^{n-1} + \cdots + c_0$. Second, we rely on the fact that a monic polynomial can be written in factored form as

$$p(z) = (z - r_1) \cdots (z - r_n) \quad (2.7)$$

based on the roots guaranteed by the fundamental theorem of algebra. Let r be a root of $p(z)$ of multiplicity m and $q(z) = z^n + d_{n-1}z^{n-1} + \cdots + d_0$ be a second monic polynomial of the same degree n as $p(z)$. We now interpret continuity to mean that for every $\epsilon > 0$, there exists a $\delta > 0$ such that $q(z)$ has at least m roots within ϵ of r whenever the coefficient vector \mathbf{d} of $q(z)$ satisfies $\|\mathbf{d} - \mathbf{c}\| < \delta$. Here we use the Euclidean norm on \mathbb{R}^{2n} . In proving this result, we need the simple bound

$$|r_j| \leq \max \left\{ 1, \sum_{i=0}^{n-1} |c_i| \right\}.$$

on the roots of a monic polynomial $p(z)$ in terms of its coefficients. The proof of the bound is an immediate consequence of the identity

$$r_j = - \sum_{i=0}^{n-1} c_i r_j^{i-n+1}.$$

We are now in a position to verify the asserted continuous dependence. Suppose it fails for the polynomial $p(z)$ and the specified root r . Then for some $\epsilon > 0$ there exists a sequence $q_k(z)$ of monic polynomials of degree n with fewer than m roots within ϵ of r but whose coefficients d_{ki} converge to the coefficients c_i . Since the coefficients of the $q_k(z)$ converge, by the above inequality, the roots s_{ki} of the $q_k(z)$ are bounded. We can therefore extract a subsequence $q_{k_l}(z)$ whose roots converge to the complex numbers t_i . At most $m - 1$ of the t_i equal r . The representation

$$p(z) = \lim_{l \rightarrow \infty} q_{k_l}(z) = (z - t_1) \cdots (z - t_n)$$

is at odds with the representation (2.7) of $p(z)$. Indeed, one has m roots equal to r , and the other has at most $m - 1$ roots equal to r . This contradiction proves the claimed continuity of the roots.

As an illustration consider the quadratic $p(z) = z^2 - 2z + 1 = (z - 1)^2$ with the root 1 of multiplicity 2. For $\delta > 0$ small the related polynomial $z^2 - 2z + 1 - \delta$ has the real roots $1 \pm \sqrt{\delta}$ while the polynomial $z^2 - 2z + 1 + \delta$ has the complex roots $1 \pm \sqrt{-\delta}$. A more important application concerns the continuity of the eigenvalues of a matrix. Suppose the sequence \mathbf{M}_k of square matrices converges to the square matrix \mathbf{M} . Then the sequence of characteristic polynomials $\det(z\mathbf{I} - \mathbf{M}_k)$ converges to the characteristic polynomial $\det(z\mathbf{I} - \mathbf{M})$. It follows that the eigenvalues of \mathbf{M}_k converge to the eigenvalues of \mathbf{M} in the sense just explained. ■

A function $f(\mathbf{x})$ is said to be uniformly continuous on its domain S if for every $\epsilon > 0$ there exists a $\delta > 0$ such that $\|f(\mathbf{y}) - f(\mathbf{x})\| < \epsilon$ whenever $\|\mathbf{y} - \mathbf{x}\| < \delta$. This sounds like ordinary continuity, but the chosen δ does not depend on the pivotal point $\mathbf{x} \in S$. One of the virtues of a compact domain is that it forces uniform continuity.

Proposition 2.5.5 (Heine) *Every continuous function $f(\mathbf{x})$ from a compact set S of \mathbb{R}^m into \mathbb{R}^n is uniformly continuous.*

Proof: Suppose $f(\mathbf{x})$ fails to be uniformly continuous. Then for some $\epsilon > 0$, there exist sequences \mathbf{x}_i and \mathbf{y}_i from S such that $\lim_{i \rightarrow \infty} \|\mathbf{x}_i - \mathbf{y}_i\| = 0$ and $\|f(\mathbf{x}_i) - f(\mathbf{y}_i)\| \geq \epsilon$. Since S is compact, we can extract a subsequence of \mathbf{x}_i that converges to a point $\mathbf{u} \in S$. Along the corresponding subsequence of \mathbf{y}_i we can extract a subsubsequence that converges to a point $\mathbf{v} \in S$. Substituting the constructed subsubsequences for \mathbf{x}_i and \mathbf{y}_i if necessary, we may assume that \mathbf{x}_i and \mathbf{y}_i both converge to the same limit $\mathbf{u} = \mathbf{v}$. The condition $\|f(\mathbf{x}_i) - f(\mathbf{y}_i)\| \geq \epsilon$ now contradicts the continuity of $f(\mathbf{x})$ at \mathbf{u} . ■

Example 2.5.9 Rigid Motions

Uniform continuity certainly appears in the absence of compactness. One spectacular example is a rigid motion. By this we mean a function $f(\mathbf{x})$ of \mathbb{R}^n into itself with the property $\|f(\mathbf{y}) - f(\mathbf{x})\| = \|\mathbf{y} - \mathbf{x}\|$ for every choice of \mathbf{x} and \mathbf{y} . We can better understand the rigid motion $f(\mathbf{x})$ by investigating the translated rigid motion $g(\mathbf{x}) = f(\mathbf{x}) - f(\mathbf{0})$ that maps the origin $\mathbf{0}$ into itself. Because $g(\mathbf{x})$ preserves distances, it also preserves inner products. This fact is evident from the equalities

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}\|^2 &= \|\mathbf{y}\|^2 - 2\mathbf{y}^* \mathbf{x} + \|\mathbf{x}\|^2 \\ \|g(\mathbf{y}) - g(\mathbf{x})\|^2 &= \|g(\mathbf{y})\|^2 - 2g(\mathbf{y})^* g(\mathbf{x}) + \|g(\mathbf{x})\|^2 \\ \|g(\mathbf{y})\|^2 &= \|\mathbf{y}\|^2 \\ \|g(\mathbf{x})\|^2 &= \|\mathbf{x}\|^2. \end{aligned}$$

The inner product identity

$$g(\mathbf{y})^* g(\mathbf{x}) = \mathbf{y}^* \mathbf{x}$$

is only possible if $g(\mathbf{y})$ is linear. To demonstrate this assertion, note that $g(\mathbf{x})$ maps the standard orthonormal basis $\mathbf{e}_1, \dots, \mathbf{e}_n$ onto the orthonormal basis $g(\mathbf{e}_1), \dots, g(\mathbf{e}_n)$. Because

$$\begin{aligned} g(\alpha \mathbf{x} + \beta \mathbf{y})^* g(\mathbf{e}_i) &= (\alpha \mathbf{x} + \beta \mathbf{y})^* \mathbf{e}_i \\ &= \alpha \mathbf{x}^* \mathbf{e}_i + \beta \mathbf{y}^* \mathbf{e}_i \\ &= [\alpha g(\mathbf{x}) + \beta g(\mathbf{y})]^* g(\mathbf{e}_i) \end{aligned}$$

holds for all i , it follows that $g(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha g(\mathbf{x}) + \beta g(\mathbf{y})$. In other words, $g(\mathbf{x})$ is linear. The linear transformations that preserve angles and distances are precisely the orthogonal transformations. Thus, the rigid motion $f(\mathbf{x})$ reduces to an orthogonal transformation $U\mathbf{x}$ followed by the translation $f(\mathbf{0})$. Conversely, it is trivial to prove that every such transformation

$$f(\mathbf{x}) = U\mathbf{x} + f(\mathbf{0})$$

is a rigid motion. ■

Example 2.5.10 Multilinear Maps

A k -linear map $M[\mathbf{u}_1, \dots, \mathbf{u}_k]$ transforms points from the k -fold Cartesian product $\mathbb{R}^m \times \mathbb{R}^m \times \dots \times \mathbb{R}^m$ into points in \mathbb{R}^n and satisfies the rules

$$\begin{aligned} M[\mathbf{u}_1, \dots, c\mathbf{u}_j, \dots, \mathbf{u}_k] &= cM[\mathbf{u}_1, \dots, \mathbf{u}_j, \dots, \mathbf{u}_k] \\ M[\mathbf{u}_1, \dots, \mathbf{u}_j + \mathbf{v}_j, \dots, \mathbf{u}_k] &= M[\mathbf{u}_1, \dots, \mathbf{u}_j, \dots, \mathbf{u}_k] \\ &\quad + M[\mathbf{u}_1, \dots, \mathbf{v}_j, \dots, \mathbf{u}_k] \end{aligned}$$

for every scalar c , index j , vector \mathbf{v}_j , and combination of vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$. For example, matrix multiplication $\mathbf{u} \mapsto A\mathbf{u}$ is 1-linear and the determinant map $U \mapsto \det U$ is k -linear on the columns \mathbf{u}_j of a $k \times k$ matrix. A k -linear map into the real line ($n = 1$) is called a k -linear form. The k -linear form

$$[\mathbf{u}_1, \dots, \mathbf{u}_k] \mapsto \prod_{j=1}^k \mathbf{v}_j^* \mathbf{u}_j \quad (2.8)$$

for any fixed combination $[\mathbf{v}_1, \dots, \mathbf{v}_k]$ of vectors is often useful in applications. A k -linear map $M[\mathbf{u}_1, \dots, \mathbf{u}_j, \dots, \mathbf{u}_k]$ is said to be symmetric if

$$M[\mathbf{u}_1, \dots, \mathbf{u}_i, \dots, \mathbf{u}_j, \dots, \mathbf{u}_k] = M[\mathbf{u}_1, \dots, \mathbf{u}_j, \dots, \mathbf{u}_i, \dots, \mathbf{u}_k]$$

for all pairs of indices i and j and antisymmetric if

$$M[\mathbf{u}_1, \dots, \mathbf{u}_i, \dots, \mathbf{u}_j, \dots, \mathbf{u}_k] = -M[\mathbf{u}_1, \dots, \mathbf{u}_j, \dots, \mathbf{u}_i, \dots, \mathbf{u}_k].$$

The determinant function is antisymmetric.

One can easily check that the collection $\mathbb{L}^k(\mathbb{R}^m, \mathbb{R}^n)$ of k -linear maps from $\mathbb{R}^m \times \mathbb{R}^m \times \dots \times \mathbb{R}^m$ to \mathbb{R}^n forms a vector space under pointwise addition and scalar multiplication. Its dimension is $m^k n$. Indeed, let $\mathbf{e}_1, \dots, \mathbf{e}_m$ be a basis for \mathbb{R}^m and $\mathbf{f}_1, \dots, \mathbf{f}_n$ be a basis for \mathbb{R}^n . If $\mathbf{u}_i = \sum_{j=1}^m c_{ij} \mathbf{e}_j$, then the expansion

$$M[\mathbf{u}_1, \dots, \mathbf{u}_k] = \sum_{j_1=1}^m \dots \sum_{j_k=1}^m \left(\prod_{i=1}^k c_{i, j_i} \right) M[\mathbf{e}_{j_1}, \dots, \mathbf{e}_{j_k}] \quad (2.9)$$

correctly suggests that the k -linear maps with

$$\mathbf{M}_{j_1, \dots, j_k, l}[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_k}] = \begin{cases} \mathbf{f}_l & i_1 = j_1, \dots, i_k = j_k \\ \mathbf{0} & \text{otherwise} \end{cases}$$

constitute a basis of $\mathbb{L}^k(\mathbb{R}^m, \mathbb{R}^n)$. For a linear form $\mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]$, it is helpful to think of the numbers $\mathbf{M}[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_k}]$ as coefficients that define the linear form just as the coefficients of a matrix define the corresponding linear transformation.

Equation (2.9) also implies that $\mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]$ is a continuous function of its arguments. Therefore, the norm

$$\|\mathbf{M}\| = \sup_{\mathbf{u}_j \neq \mathbf{0} \forall j} \frac{\|\mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]\|}{\|\mathbf{u}_1\| \cdots \|\mathbf{u}_k\|} = \sup_{\|\mathbf{u}_j\|=1 \forall j} \|\mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]\|$$

on $\mathbb{L}^k(\mathbb{R}^m, \mathbb{R}^n)$ induced by the Euclidean norms on \mathbb{R}^m and \mathbb{R}^n is finite. For example, the norm of the k -linear form (2.8) is $\prod_{j=1}^k \|\mathbf{v}_j\|$. This value is attained by choosing $\mathbf{u}_j = \|\mathbf{v}_j\|^{-1} \mathbf{v}_j$ and serves as an absolute upper bound on the k -linear form on unit vectors by virtue of the Cauchy-Schwarz inequality. The inequality

$$\|\mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]\| \leq \|\mathbf{M}\| \|\mathbf{u}_1\| \cdots \|\mathbf{u}_k\| \quad (2.10)$$

is an immediate consequence of the definition of $\|\mathbf{M}\|$. Problem 33 asks the reader to verify that the map $(\mathbf{M}, \mathbf{u}_1, \dots, \mathbf{u}_k) \mapsto \mathbf{M}[\mathbf{u}_1, \dots, \mathbf{u}_k]$ is jointly continuous in its $k+1$ variables. ■

2.6 Semicontinuity

For real-valued functions, the notions of lower and upper semicontinuity are often useful substitutes for continuity. A real-valued function $f(\mathbf{x})$ with domain $T \subset \mathbb{R}^m$ is lower semicontinuous if the set $\{\mathbf{x} \in T : f(\mathbf{x}) \leq c\}$ is closed in T for every constant c . Given the duality of closed and open sets, an equivalent condition is that $\{\mathbf{x} \in T : f(\mathbf{x}) > c\}$ is open in T for every constant c . A real-valued function $g(\mathbf{x})$ is said to be upper semicontinuous if and only if $f(\mathbf{x}) = -g(\mathbf{x})$ is lower semicontinuous. Owing to this simple relationship, we will confine our attention to lower semicontinuous functions. The next proposition gives two alternative definitions.

Proposition 2.6.1 *A necessary and sufficient condition for $f(\mathbf{x})$ to be lower semicontinuous is that*

$$f(\mathbf{x}) \leq \liminf_{n \rightarrow \infty} f(\mathbf{x}_n) \quad (2.11)$$

whenever $\lim_{n \rightarrow \infty} \mathbf{x}_n = \mathbf{x}$ in T . Another necessary and sufficient condition is that the epigraph $\{(\mathbf{x}, y) \in T \times \mathbb{R} : f(\mathbf{x}) \leq y\}$ is a closed set.

Proof: Suppose $f(\mathbf{x})$ is lower semicontinuous and $\lim_{n \rightarrow \infty} \mathbf{x}_n = \mathbf{x}$. For any $\epsilon > 0$, the point \mathbf{x} lies in the open set $\{\mathbf{y} \in T : f(\mathbf{y}) > f(\mathbf{x}) - \epsilon\}$. Hence, $f(\mathbf{x}_n) > f(\mathbf{x}) - \epsilon$ for all sufficient large n . But this implies inequality (2.11). Similarly, if \mathbf{x}_n converges to \mathbf{x} and $y_n \geq f(\mathbf{x}_n)$ converges to y , then the inequality $y \geq \liminf_{n \rightarrow \infty} f(\mathbf{x}_n) \geq f(\mathbf{x})$ follows, and the epigraph is closed. Thus, both stated conditions are necessary.

For sufficiency, suppose inequality (2.11) holds. Consider a sequence \mathbf{x}_n in the set $\{\mathbf{y} \in T : f(\mathbf{y}) \leq c\}$ with limit \mathbf{x} . It is clear that $f(\mathbf{x}) \leq c$ as well. Thus, $\{\mathbf{y} \in T : f(\mathbf{y}) \leq c\}$ is closed. To deal with the second sufficient condition, suppose the epigraph is closed, but $f(\mathbf{x})$ is not lower semicontinuous. Then there exists a sequence \mathbf{x}_n converging to \mathbf{x} in T and an $\epsilon > 0$ such that $f(\mathbf{x}) - \epsilon > \liminf_{n \rightarrow \infty} f(\mathbf{x}_n)$. It follows that the pair $(\mathbf{x}_n, f(\mathbf{x}) - \epsilon)$ is in the epigraph for infinitely many n . Because the epigraph is closed, this forces the contradiction that $(\mathbf{x}, f(\mathbf{x}) - \epsilon)$ belongs to the epigraph. ■

Part of the motivation for defining semicontinuity is to generalize Proposition 2.5.4. The result stated there for global maxima holds for upper semicontinuous functions, and the result for global minima holds for lower semicontinuous functions. The proof carries over almost word for word. It is also obvious that any continuous function is lower semicontinuous, and any function that is both lower and upper semicontinuous is continuous. Fortunately, the closure properties of lower semicontinuous functions are quite flexible.

Proposition 2.6.2 *The collection of lower semicontinuous functions with common domain $T \subset \mathbb{R}^m$ satisfies the following rules:*

- (a) *If $f_k(\mathbf{x})$ is a family of lower semicontinuous functions, then $\sup_k f_k(\mathbf{x})$ is lower semicontinuous.*
- (b) *If $f_k(\mathbf{x})$ is a finite family of lower semicontinuous functions, then $\min_k f_k(\mathbf{x})$ is lower semicontinuous.*
- (c) *If $f(\mathbf{x})$ and $g(\mathbf{x})$ are lower semicontinuous, then $f(\mathbf{x}) + g(\mathbf{x})$ is lower semicontinuous.*
- (d) *If $f(\mathbf{x})$ and $g(\mathbf{x})$ are both positive and lower semicontinuous, then $f(\mathbf{x})g(\mathbf{x})$ is lower semicontinuous.*
- (e) *If $f(\mathbf{x})$ is lower semicontinuous and $g(\mathbf{x})$ is continuous with range U contained in T , then $f \circ g(\mathbf{x})$ is lower semicontinuous.*

Proof: These rules follow from the set identities

$$\begin{aligned} \{\mathbf{x} : \sup_k f_k(\mathbf{x}) > c\} &= \cup_k \{\mathbf{x} : f_k(\mathbf{x}) > c\} \\ \{\mathbf{x} : \min_k f_k(\mathbf{x}) > c\} &= \cap_k \{\mathbf{x} : f_k(\mathbf{x}) > c\} \end{aligned}$$

$$\begin{aligned} \{\mathbf{x} : f(\mathbf{x}) + g(\mathbf{x}) > c\} &= \cup_d(\{\mathbf{x} : f(\mathbf{x}) > c - d\} \cap \{\mathbf{x} : g(\mathbf{x}) > d\}) \\ \{\mathbf{x} : f(\mathbf{x})g(\mathbf{x}) > c\} &= \cup_{d>0}(\{\mathbf{x} : f(\mathbf{x}) > d^{-1}c\} \cap \{\mathbf{x} : g(\mathbf{x}) > d\}) \\ \{\mathbf{y} : f \circ g(\mathbf{y}) > c\} &= g^{-1}[\{\mathbf{x} : f(\mathbf{x}) > c\}] \end{aligned}$$

and the properties of open sets and continuous functions summarized in Propositions 2.4.2 and 2.5.2. ■

Example 2.6.1 Row and Column Rank

Every $m \times n$ matrix \mathbf{A} has a well defined nullity and rank. Although these are not continuous functions of \mathbf{A} , the former function is upper semicontinuous, and the latter function is lower semicontinuous. In view of the dimension identity $\text{nullity}(\mathbf{A}) = n - \text{rank}(\mathbf{A})$, to validate both claims it suffices to show that $\text{rank}(\mathbf{A})$ is lower semicontinuous. Consider an arbitrary constant c and an arbitrary matrix $\mathbf{A} = (a_{ij})$ with $\text{rank}(\mathbf{A}) > c$. If we abbreviate $\text{rank}(\mathbf{A}) = r$, then there exist row indices $1 \leq i_1 < \cdots < i_r \leq m$ and column indices $1 \leq j_1 < \cdots < j_r \leq n$ such that the submatrix

$$\begin{pmatrix} a_{i_1 j_1} & \cdots & a_{i_1 j_r} \\ \vdots & \ddots & \vdots \\ a_{i_r j_1} & \cdots & a_{i_r j_r} \end{pmatrix}$$

has nonzero determinant. Because the determinant function is continuous, the same submatrix has nonvanishing determinant for all $m \times n$ matrices \mathbf{B} close to \mathbf{A} . It follows that $\{\mathbf{A} : \text{rank}(\mathbf{A}) > c\}$ is an open set and therefore that $\text{rank}(\mathbf{A})$ is lower semicontinuous. ■

2.7 Connectedness

Roughly speaking, a set is disconnected if it can be split into two pieces sharing no boundary. A set is connected if it is not disconnected. One way of making this vague distinction precise is to consider a set S disconnected if there exists a real-valued continuous function $\phi(\mathbf{x})$ defined on S and having range $\{0, 1\}$. The nonempty subsets $A = \phi^{-1}(0)$ and $B = \phi^{-1}(1)$ then constitute the two disconnected pieces of S . According to part (b) of Proposition 2.5.2, both A and B are closed. Because one is the complement of the other, both are also open.

Arcwise connectedness is a variation on the theme of connectedness. A set is said to be arcwise connected if for any pair of points \mathbf{x} and \mathbf{y} of the set there is a continuous function $f(t)$ from the interval $[0, 1]$ into the set satisfying $f(0) = \mathbf{x}$ and $f(1) = \mathbf{y}$. We will see shortly that arcwise connectedness implies connectedness. On open sets, the two notions coincide.

Can we identify the connected subsets of the real line? Intuition suggests that the only connected subsets are intervals. Here a single point x is viewed

as the interval $[x, x]$. Suppose S is a connected subset of \mathbb{R} , and let a and b be two points of S . In order for S to be an interval, every point $c \in (a, b)$ should be in S . If S fails to contain an intermediate point c , then we can define a continuous function $\phi(x)$ disconnecting S by taking $\phi(x) = 0$ for $x < c$ and $\phi(x) = 1$ for $x > c$. Thus, every connected subset must be an interval.

To prove the converse, suppose a disconnecting function $\phi(x)$ lives on an interval. Select points a and b of the interval with $\phi(a) = 0$ and $\phi(b) = 1$. Without loss of generality we can take $a < b$. On $[a, b]$ we now carry out the bisection strategy of Example 2.3.1, selecting the right or left subinterval at each stage so that the values of $\phi(x)$ at the endpoints of the selected subinterval disagree. Eventually, bisection leads to a subinterval contradicting the uniform continuity of $\phi(x)$ on $[a, b]$. Indeed, there is a number δ such that $|\phi(y) - \phi(x)| < 1$ whenever $|y - x| < \delta$; at some stage, the length of the subinterval containing points with both values of $\phi(x)$ falls below δ .

This result is the first of four characterizing connected sets.

Proposition 2.7.1 *Connected subsets of \mathbb{R}^n have the following properties:*

- (a) *A subset of the real line is connected if and only if it is an interval.*
- (b) *The image of a connected set under a continuous function is connected.*
- (c) *The union $S = \cup_{\alpha} S_{\alpha}$ of an arbitrary collection of connected subsets is connected if one of the sets S_{β} has a nonempty intersection $S_{\beta} \cap S_{\alpha}$ with every other set S_{α} .*
- (d) *Every arcwise connected set S is connected.*

Proof: To prove part (b) let $f(\mathbf{x})$ be a continuous map from a connected set $S \subset \mathbb{R}^m$ into \mathbb{R}^n . If the image $f(S)$ is disconnected, then there is a continuous function $\phi(\mathbf{x})$ disconnecting it. The composition $\phi \circ f(\mathbf{x})$ is continuous by part (g) of Proposition 2.5.1 and serves to disconnect S , contradicting the connectedness of S . To prove (c) suppose that the continuous function $\phi(\mathbf{x})$ disconnects the union S . Then there exists $\mathbf{y} \in S_{\alpha_1}$ and $\mathbf{z} \in S_{\alpha_2}$ with $\phi(\mathbf{y}) = 0$ and $\phi(\mathbf{z}) = 1$. Choose $\mathbf{u} \in S_{\beta} \cap S_{\alpha_1}$ and $\mathbf{v} \in S_{\beta} \cap S_{\alpha_2}$. If $\phi(\mathbf{u}) \neq \phi(\mathbf{v})$, then $\phi(\mathbf{x})$ disconnects S_{β} . If $\phi(\mathbf{u}) = \phi(\mathbf{v})$, then $\phi(\mathbf{y}) \neq \phi(\mathbf{u})$ or $\phi(\mathbf{z}) \neq \phi(\mathbf{v})$. In the former case $\phi(\mathbf{x})$ disconnects S_{α_1} , and in the latter case $\phi(\mathbf{x})$ disconnects S_{α_2} . Finally, to prove part (d), suppose the arcwise connected set S fails to be connected. Then there exists a continuous disconnecting function $\phi(\mathbf{x})$ with $\phi(\mathbf{y}) = 0$ and $\phi(\mathbf{z}) = 1$. Let $f(t)$ be an arc in S connecting \mathbf{y} and \mathbf{z} . The continuous function $\phi \circ f(t)$ then serves to disconnect $[0, 1]$. ■

Example 2.7.1 *The Intermediate Value Property*

Consider a continuous function $f(x)$ from an interval $[a, b]$ to the real line. The intermediate value theorem asserts that the image $f([a, b])$ coincides with the interval $[\min f(x), \max f(x)]$. This theorem, which is a

consequence of properties (a) and (b) of Proposition 2.7.1, has many applications. For example, suppose $g(x)$ is a continuous function from $[0, 1]$ into $[0, 1]$. If $f(x) = g(x) - x$, then it is obvious that $f(0) \geq 0$ and $f(1) \leq 0$. It follows that $f(x) = 0$ for some x . In other words, $g(x)$ has a fixed point satisfying $g(x) = x$. ■

Example 2.7.2 Connectedness of Spheres

The set $S(\mathbf{x}, r)$ in \mathbb{R}^n is the image of the continuous map $\mathbf{y} \mapsto \mathbf{x} + r\mathbf{y}/\|\mathbf{y}\|$ of the domain $T = \mathbb{R}^n \setminus \mathbf{0}$. Hence, to prove connectedness when $n > 1$, it suffices to prove that T is connected. To achieve this, we argue that T is arcwise connected. Consider two points \mathbf{u} and \mathbf{v} in T . If $\mathbf{0}$ does not lie on the line segment between \mathbf{u} and \mathbf{v} , then we can use the function $f(t) = \mathbf{u} + t(\mathbf{v} - \mathbf{u})$ to connect \mathbf{u} and \mathbf{v} . If $\mathbf{0}$ lies on the line segment, choose any \mathbf{w} not on the line determined by \mathbf{u} and \mathbf{v} . Now the continuous function

$$f(t) = \begin{cases} \mathbf{u} + 2t(\mathbf{w} - \mathbf{u}) & t \in [0, \frac{1}{2}] \\ \mathbf{w} + (2t - 1)(\mathbf{v} - \mathbf{w}) & t \in [\frac{1}{2}, 1] \end{cases}$$

connects \mathbf{u} and \mathbf{v} . The sphere $S(x, r)$ in \mathbb{R} reduces to the two points $x - r$ and $x + r$ and is disconnected. ■

2.8 Uniform Convergence

Many delicate issues of analysis revolve around the question of whether a given property of a sequence of functions $f_m(\mathbf{x})$ is preserved under a passage to a limit. As a simple example, consider the sequence $f_m(x) = x^m$ of continuous functions defined on the unit interval $[0, 1]$. It is clear that $f_m(x)$ converges pointwise to the discontinuous function

$$f(x) = \begin{cases} 0 & 0 \leq x < 1 \\ 1 & x = 1. \end{cases}$$

The failure of $f(x)$ to be continuous suggests that an additional hypothesis must be imposed. The key hypothesis is uniform convergence. This requires for each $\epsilon > 0$ that there exists an integer k such that $|f_m(\mathbf{x}) - f(\mathbf{x})| < \epsilon$ for all $m \geq k$ and all \mathbf{x} . Here the adjective “uniform” refers to the assumption that the same k works for all \mathbf{x} . Of course, k is allowed to depend on ϵ .

Proposition 2.8.1 *Suppose the sequence of continuous functions $f_m(\mathbf{x})$ maps a domain $D \subset \mathbb{R}^p$ into \mathbb{R}^q . If $f_m(\mathbf{x})$ converges uniformly to $f(\mathbf{x})$ on D , then $f(\mathbf{x})$ is also continuous.*

Proof: Choose $\mathbf{y} \in D$ and $\epsilon > 0$, and take k so that $\|f_m(\mathbf{x}) - f(\mathbf{x})\| < \frac{\epsilon}{3}$ for all $m \geq k$ and \mathbf{x} . By virtue of the continuity of $f_k(\mathbf{x})$, there is a $\delta > 0$

such that $\|f_k(\mathbf{x}) - f_k(\mathbf{y})\| < \frac{\epsilon}{3}$ whenever $\|\mathbf{x} - \mathbf{y}\| < \delta$. Assuming that \mathbf{y} is fixed and $\|\mathbf{x} - \mathbf{y}\| < \delta$, we have

$$\begin{aligned} \|f(\mathbf{x}) - f(\mathbf{y})\| &\leq \|f(\mathbf{x}) - f_k(\mathbf{x})\| + \|f_k(\mathbf{x}) - f_k(\mathbf{y})\| + \|f_k(\mathbf{y}) - f(\mathbf{y})\| \\ &< \frac{\epsilon}{3} + \frac{\epsilon}{3} + \frac{\epsilon}{3} \\ &= \epsilon. \end{aligned}$$

This shows that $f(\mathbf{x})$ is continuous at \mathbf{y} . ■

Example 2.8.1 Weierstrass M -Test

Suppose the entries $g_k(\mathbf{x})$ of a sequence of continuous functions satisfy $\|g_k(\mathbf{x})\| \leq M_k$, where $\sum_{k=1}^{\infty} M_k < \infty$. Then Cauchy's criterion and Proposition 2.8.1 together imply that the partial sums $f_l(\mathbf{x}) = \sum_{k=1}^l g_k(\mathbf{x})$ converge uniformly to the continuous function $f(\mathbf{x}) = \sum_{k=1}^{\infty} g_k(\mathbf{x})$. ■

2.9 Problems

- Let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be points in \mathbb{R}^n . State and prove a necessary and sufficient condition under which the Euclidean norm equality

$$\|\mathbf{x}_1 + \dots + \mathbf{x}_m\| = \|\mathbf{x}_1\| + \dots + \|\mathbf{x}_m\|$$

holds. (Hints: Square and expand both sides. Use the necessary and sufficient conditions of the Cauchy-Schwarz inequality term by term.)

- Show that it is possible to choose $n + 1$ points $\mathbf{x}_0, \mathbf{x}_1, \dots, \mathbf{x}_n$ in \mathbb{R}^n such that $\|\mathbf{x}_i\| = 1$ for all i and $\|\mathbf{x}_i - \mathbf{x}_j\| = \|\mathbf{x}_k - \mathbf{x}_l\|$ for all pairs $i \neq j$ and $k \neq l$. These points define a regular simplex with vertices on the unit sphere. (Hint: One possibility is to take $\mathbf{x}_0 = n^{-1/2}\mathbf{1}$ and $\mathbf{x}_i = a\mathbf{1} + b\mathbf{e}_i$ for $i \geq 1$, where

$$a = -\frac{1 + \sqrt{n+1}}{n^{3/2}}, \quad b = \sqrt{\frac{n+1}{n}}.$$

Any rotated version of these points also works.)

- Show that

$$\|\mathbf{x}\|_q \leq \|\mathbf{x}\|_p \tag{2.12}$$

$$\|\mathbf{x}\|_p \leq n^{\frac{1}{p} - \frac{1}{q}} \|\mathbf{x}\|_q \tag{2.13}$$

when p and q are chosen from $\{1, 2, \infty\}$ and $p < q$. Here $\|\mathbf{x}\|_2$ is the Euclidean norm on \mathbb{R}^n . These inequalities are sharp. Equality holds in inequality (2.12) when $\mathbf{x} = (1, 0, \dots, 0)^*$, and equality holds in inequality (2.13) when $\mathbf{x} = (1, 1, \dots, 1)^*$.

4. Show that $\|\mathbf{x}\|^2 \leq \|\mathbf{x}\|_\infty \|\mathbf{x}\|_1 \leq \sqrt{n} \|\mathbf{x}\|^2$ for any vector $\mathbf{x} \in \mathbb{R}^n$.
5. Prove that $1 \leq \|\mathbf{I}\|_\dagger$ and $\|\mathbf{M}\|_\dagger^{-1} \leq \|\mathbf{M}^{-1}\|_\dagger$ for any matrix norm on square matrices satisfying the defining properties (a) through (e) of Sect. 2.2.
6. Set $\|\mathbf{M}\|_{\max} = \max_{i,j} |m_{ij}|$ for $\mathbf{M} = (m_{ij})$. Show that this defines a vector norm but not a matrix norm on $n \times n$ matrices \mathbf{M} .
7. Let \mathbf{M} be an $m \times n$ matrix. Prove that its spectral norm satisfies

$$\|\mathbf{M}\| = \sup_{\|\mathbf{v}\|=1} \|\mathbf{M}\mathbf{v}\| = \sup_{\|\mathbf{u}\|=1, \|\mathbf{v}\|=1} \mathbf{u}^* \mathbf{M} \mathbf{v}$$

for $\mathbf{u} \in \mathbb{R}^m$ and $\mathbf{v} \in \mathbb{R}^n$.

8. Demonstrate that the spectral norm on $m \times n$ matrices satisfies $\|\mathbf{U}\mathbf{M}\| = \|\mathbf{M}\| = \|\mathbf{M}\mathbf{V}\|$ for all orthogonal matrices \mathbf{U} and \mathbf{V} of the right dimensions. Show that the Frobenius norm satisfies the same orthogonal invariance principle.
9. Show that an $m \times n$ matrix $\mathbf{M} = (m_{ij})$ has the matrix norms

$$\begin{aligned} \|\mathbf{M}\|_1 &= \max_{1 \leq j \leq n} \sum_{i=1}^m |m_{ij}| \\ \|\mathbf{M}\|_\infty &= \max_{1 \leq i \leq m} \sum_{j=1}^n |m_{ij}| \end{aligned}$$

induced by the vector norms $\|\mathbf{x}\|_1$ and $\|\mathbf{x}\|_\infty$.

10. Let \mathbf{M} be an $m \times n$ matrix of full column rank n . Prove that there exists a positive constant c such that $\|\mathbf{M}\mathbf{y}\| \geq c\|\mathbf{y}\|$ for all $\mathbf{y} \in \mathbb{R}^n$.
11. Demonstrate properties (2.4) and (2.5) of the limit superior and limit inferior. Also check that the sequence x_n has a limit if and only if equality holds in inequality (2.5).
12. Let $l = \limsup_{n \rightarrow \infty} x_n$. Show that:
 - (a) $l = -\infty$ if and only if $\lim_{n \rightarrow \infty} x_n = -\infty$.
 - (b) $l = +\infty$ if and only if for every positive integer m and real r there exists an $n \geq m$ with $x_n > r$.
 - (c) l is finite if and only if (a) for every $\epsilon > 0$ there is an m such that $n \geq m$ implies $x_n < l + \epsilon$ and (b) for every $\epsilon > 0$ and every m there is an $n \geq m$ such that $x_n > l - \epsilon$.

Similar properties hold for the limit inferior.

13. For any sequence of real numbers x_n , prove that

$$\inf_n x_n \leq \liminf_{n \rightarrow \infty} x_n \text{ and } \limsup_{n \rightarrow \infty} x_n \leq \sup_n x_n.$$

If y_n is a second sequence of real numbers, then prove that

$$\begin{aligned} \liminf_{n \rightarrow \infty} (x_n + y_n) &\geq \liminf_{n \rightarrow \infty} x_n + \liminf_{n \rightarrow \infty} y_n \\ \limsup_{n \rightarrow \infty} (x_n + y_n) &\leq \limsup_{n \rightarrow \infty} x_n + \limsup_{n \rightarrow \infty} y_n. \end{aligned}$$

Finally, if $x_n \leq y_n$ for all n , then prove that

$$\liminf_{n \rightarrow \infty} x_n \leq \liminf_{n \rightarrow \infty} y_n \text{ and } \limsup_{n \rightarrow \infty} x_n \leq \limsup_{n \rightarrow \infty} y_n.$$

14. Let x_n be a sequence of nonnegative real numbers with

$$x_{n+1} \leq x_n + \frac{1}{n^2}$$

for all $n \geq 1$. Show that $\lim_{n \rightarrow \infty} x_n$ exists [69].

15. Let \mathbf{x}_m be a convergent sequence in \mathbf{R}^n with limit \mathbf{x} . Prove that the sequence $\mathbf{s}_m = (\mathbf{x}_1 + \cdots + \mathbf{x}_m)/m$ of arithmetic means converges to \mathbf{x} .

16. Show that

$$\lim_{x \rightarrow \infty} p(x)e^{-x} = 0$$

for every polynomial $p(x)$.

17. Prove that the set of invertible matrices is open and that the sets of symmetric and orthogonal matrices are closed in \mathbf{R}^{n^2} .

18. A square matrix is nilpotent if $\mathbf{A}^k = \mathbf{0}$ for some positive integer k . If \mathbf{A} and \mathbf{B} are nilpotent, then show that $\mathbf{A} + \mathbf{B}$ need not be nilpotent. If we add the hypothesis that \mathbf{A} and \mathbf{B} commute, then show that $\mathbf{A} + \mathbf{B}$ is nilpotent. Use Example 2.3.3 to construct the inverses of the matrices $\mathbf{I} + \mathbf{A}$ and $\mathbf{I} - \mathbf{A}$ for \mathbf{A} nilpotent [69].

19. Show that $e^{-\mathbf{M}}$ is the matrix inverse of $e^{\mathbf{M}}$. A skew symmetric matrix \mathbf{M} satisfies $\mathbf{M}^* = -\mathbf{M}$. Show that $e^{\mathbf{M}}$ is orthogonal when \mathbf{M} is skew symmetric.

20. Demonstrate that the matrix exponential function $\mathbf{M} \mapsto e^{\mathbf{M}}$ is continuous. (Hint: Apply the Weierstrass M -test.)

21. Demonstrate that the function

$$f(\mathbf{x}) = \begin{cases} \frac{x_1 x_2}{x_1^2 - x_2^2} & |x_1| \neq |x_2| \\ 0 & \text{otherwise} \end{cases}$$

on \mathbf{R}^2 is discontinuous at $\mathbf{0}$.

22. Let $f(\mathbf{x})$ and $g(\mathbf{x})$ be real-valued continuous functions defined on the same domain. Prove that $\max\{f(\mathbf{x}), g(\mathbf{x})\}$ and $\min\{f(\mathbf{x}), g(\mathbf{x})\}$ are continuous functions. Prove that the function $f(\mathbf{x}) = \max_i x_i$ is continuous on \mathbb{R}^n .

23. Define the function

$$f(x) = \begin{cases} x & x \text{ is rational} \\ 1 - x & x \text{ is irrational} \end{cases}$$

on $[0, 1]$. At what points is $f(x)$ continuous? What is the image $f([0, 1])$?

24. Give an example of a continuous function that does not map an open set to an open set. Give another example of a continuous function that does not map a closed set to a closed set.

25. Show that the set of $n \times n$ orthogonal matrices is compact. (Hint: Show that every orthogonal matrix \mathbf{O} has norm $\|\mathbf{O}\| = 1$.)

26. Let $f(\mathbf{x})$ be a continuous function from a compact set $S \subset \mathbb{R}^m$ into \mathbb{R}^n . If $f(\mathbf{x})$ is one-to-one, then demonstrate that the inverse function $f^{-1}(\mathbf{y})$ is continuous from $f(S)$ to S .

27. Let $C = A \times B$ be the Cartesian product of two subsets $A \subset \mathbb{R}^m$ and $B \subset \mathbb{R}^n$. Prove that:

- (a) C is closed in \mathbb{R}^{m+n} if both A and B are closed.
- (b) C is open in \mathbb{R}^{m+n} if both A and B are open.
- (c) C is compact in \mathbb{R}^{m+n} if both A and B are compact.
- (d) C is connected in \mathbb{R}^{m+n} if both A and B are connected.

28. Prove the converse of each of the assertions in Problem 27.

29. Without appeal to Proposition 2.5.5, show that every polynomial on \mathbb{R} is uniformly continuous on a compact interval $[a, b]$.

30. Let $f(x)$ be uniformly continuous on \mathbb{R} and satisfy $f(0) = 0$. Demonstrate that there exists a nonnegative constant c such that

$$|f(x)| \leq 1 + c|x|$$

for all x [69].

31. Suppose that $f(x)$ is continuous on $[0, \infty)$ and $\lim_{x \rightarrow \infty} f(x)$ exists and is finite. Prove that $f(x)$ is uniformly continuous on $[0, \infty)$ [69].

32. Characterize those maps $f(\mathbf{x})$ from \mathbb{R}^n into itself that have the property $\|f(\mathbf{y}) - f(\mathbf{x})\| = c\|\mathbf{y} - \mathbf{x}\|$ for all \mathbf{x} and \mathbf{y} . Here the constant c need not equal 1.

33. Prove that the multilinear map $(M, \mathbf{u}_1, \dots, \mathbf{u}_k) \mapsto M[\mathbf{u}_1, \dots, \mathbf{u}_k]$ is jointly continuous in its $k + 1$ variables. (Hint: Write

$$\begin{aligned} M[\mathbf{u}_1, \dots, \mathbf{u}_k] - N[\mathbf{v}_1, \dots, \mathbf{v}_k] &= (M - N)[\mathbf{u}_1, \dots, \mathbf{u}_k] \\ &+ N[\mathbf{u}_1 - \mathbf{v}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \\ &+ N[\mathbf{v}_1, \mathbf{u}_2 - \mathbf{v}_2, \dots, \mathbf{u}_k] \\ &\quad \vdots \\ &+ N[\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{u}_k - \mathbf{v}_k] \end{aligned}$$

and take norms.)

34. Let $M[\mathbf{u}_1, \dots, \mathbf{u}_k]$ be a symmetric k -linear map. Demonstrate that

$$M[\mathbf{u}_1, \dots, \mathbf{u}_k] = \frac{1}{2^k k!} \sum \epsilon_1 \cdots \epsilon_k M[(\epsilon_1 \mathbf{u}_1 + \cdots + \epsilon_k \mathbf{u}_k)^k],$$

where the sum ranges over all combinations of $\epsilon_1 = \pm 1, \dots, \epsilon_k = \pm 1$. Hence, a symmetric k -linear map is determined by its values on the diagonal of its domain.

35. Continuing Problem 34, define the alternative norm

$$\|M\|_{\text{sym}} = \sup_{\mathbf{u} \neq \mathbf{0}} \frac{\|M[\mathbf{u}^k]\|}{\|\mathbf{u}\|^k} = \sup_{\|\mathbf{u}\|=1} \|M[\mathbf{u}^k]\|.$$

Prove the inequalities

$$\|M\|_{\text{sym}} \leq \|M\| \leq \frac{k^k}{k!} \|M\|_{\text{sym}}.$$

36. Show that the indicator function of an open set is lower semicontinuous and that the indicator function of a closed set is upper semicontinuous. Also show that floor function $f(x) = \lfloor x \rfloor$ is upper semicontinuous and that the ceiling function $f(x) = \lceil x \rceil$ is lower semicontinuous.
37. Suppose the n numbers x_1, \dots, x_n lie on $[0, 1]$. Prove that the function

$$f(x) = \frac{1}{n} \sum_{i=1}^n |x - x_i|$$

attains the value $\frac{1}{2}$ for some $x \in [0, 1]$. (Hint: Consider $f(0)$ and $f(1)$.)

38. Show that a hyperplane $\{\mathbf{x} : \mathbf{z}^* \mathbf{x} = c\}$ in \mathbb{R}^n is connected but that its complement is disconnected.

39. If the real-valued function $f(x)$ on $[a, b]$ is continuous and one-to-one, then prove that $f(x)$ is either strictly increasing or strictly decreasing.
40. Demonstrate that a polynomial of odd degree possesses at least one real root.
41. Prove that the closure of a connected set is connected.
42. Suppose T is a connected set in \mathbb{R}^n . Define

$$U_\epsilon = \{\mathbf{y} \in \mathbb{R}^n : \text{dist}(\mathbf{y}, T) < \epsilon\}$$

$$V_\epsilon = \{\mathbf{y} \in \mathbb{R}^n : \text{dist}(\mathbf{y}, T) \leq \epsilon\}$$

for $\epsilon > 0$ and $\text{dist}(\mathbf{y}, T) = \inf_{\mathbf{x} \in T} \|\mathbf{y} - \mathbf{x}\|$. Demonstrate that U_ϵ and V_ϵ are connected. (Hints: V_ϵ is the closure of U_ϵ . For U_ϵ argue by contradiction using the definition of a connected set.)

43. On what domains do the sequences of functions

- (a) $f_n(x) = (1 + x/n)^n$
 (b) $f_n(x) = nx/(1 + n^2x^2)$
 (c) $f_n(x) = n^x$
 (d) $f_n(x) = x^{-1} \sin(nx)$
 (e) $f_n(x) = xe^{-nx}$
 (f) $f_n(x) = x^{2n}/(1 + x^{2n})$

converge [68]? On what domains do they converge uniformly?

44. Suppose that $f(x)$ is a function from the real line to itself satisfying $f(x + y) = f(x) + f(y)$ for all x and y . If $f(x)$ is continuous at a single point, then show that $f(x) = cx$ for some constant c . (Hints: Prove that $f(x)$ is continuous everywhere and that $f(q) = f(1)q$ for all rational numbers q .)
45. Suppose that $g(x)$ is a function from the real line to itself satisfying $g(x + y) = g(x)g(y)$ for all x and y . If $g(x)$ is continuous at a single point, then prove that either $g(x)$ is identically 0 or that there exists a positive constant d with $g(x) = d^x$. (Hint: Show that either $g(x)$ is identically 0 or that $g(x)$ is positive for all x . In the latter case, take logarithms and reduce to the previous problem.)
46. Suppose the real-valued function $f(\mathbf{x}, \mathbf{y})$ is jointly continuous in its two vector arguments and C is a compact set. Show that the functions

$$g(\mathbf{x}) = \inf_{\mathbf{y} \in C} f(\mathbf{x}, \mathbf{y}) \quad \text{and} \quad h(\mathbf{x}) = \sup_{\mathbf{y} \in C} f(\mathbf{x}, \mathbf{y})$$

are continuous.

3

The Gauge Integral

3.1 Introduction

Much of calculus deals with the interplay between differentiation and integration. The antiquated term “antidifferentiation” emphasizes the fact that differentiation and integration are inverses of one another. We will take it for granted that readers are acquainted with the mechanics of integration. The current chapter develops just enough integration theory to make our development of differentiation in Chap. 4 and the calculus of variations in Chap. 17 respectable. It is only fair to warn readers that in other chapters a few applications to probability and statistics will assume familiarity with properties of the expectation operator not covered here.

The first successful effort to put integration on a rigorous basis was undertaken by Riemann. In the early twentieth century, Lebesgue defined a more sophisticated integral that addresses many of the limitations of the Riemann integral. However, even Lebesgue’s integral has its defects. In the past few decades, mathematicians such as Henstock and Kurzweil have expanded the definition of integration on the real line to include a wider variety of functions. The new integral emerging from these investigations is called the gauge integral or generalized Riemann integral [7, 68, 108, 193, 250, 255, 278]. The gauge integral subsumes the Riemann integral, the Lebesgue integral, and the improper integrals met in traditional advanced calculus courses. In contrast to the Lebesgue integral, the integrands of the gauge integral are not necessarily absolutely integrable.

It would take us too far afield to develop the gauge integral in full generality. Here we will rest content with proving some of its elementary properties. One of the advantages of the gauge integral is that many theorems hold with fewer qualifications. The fundamental theorem of calculus is a case in point. The commonly stated version of the fundamental theorem concerns a differentiable function $f(x)$ on an interval $[a, b]$. As all students of calculus know,

$$\int_a^b f'(x) dx = f(b) - f(a).$$

Although this version is true for the gauge integral, it does not hold for the Lebesgue integral because the mere fact that $f'(x)$ exists throughout $[a, b]$ does not guarantee that it is Lebesgue integrable.

This quick description of the gauge integral is not intended to imply that the gauge integral is uniformly superior to the Lebesgue integral and its extensions. Certainly, probability theory would be severely handicapped without the full flexibility of modern measure theory. Furthermore, the advanced theory of the gauge integral is every bit as difficult as the advanced theory of the Lebesgue integral. For pedagogical purposes, however, one can argue that a student's first exposure to the theory of integration should feature the gauge integral. As we shall see, many of the basic properties of the gauge integral flow directly from its definition. As an added dividend, gauge functions provide an alternative approach to some of the material of Chap. 2.

3.2 Gauge Functions and δ -Fine Partitions

The gauge integral is defined through gauge functions. A gauge function is nothing more than a positive function $\delta(t)$ defined on a finite interval $[a, b]$. In approximating the integral of a function $f(t)$ over $[a, b]$ by a finite Riemann sum, it is important to sample the function most heavily in those regions where it changes most rapidly. Now by a Riemann sum we mean a sum

$$S(f, \pi) = \sum_{i=0}^{n-1} f(t_i)(s_{i+1} - s_i),$$

where the mesh points $a = s_0 < s_1 < \cdots < s_n = b$ form a partition π of $[a, b]$, and the tags t_i are chosen so that $t_i \in [s_i, s_{i+1}]$. If $\delta(t_i)$ measures the rapidity of change of $f(t)$ near t_i , then it makes sense to take $\delta(t)$ small in regions of rapid change and to force s_i and s_{i+1} to belong to the interval $(t_i - \delta(t_i), t_i + \delta(t_i))$. A tagged partition with this property is called a δ -fine partition. Our first proposition relieves our worry that δ -fine partitions exist.

Proposition 3.2.1 (Cousin's Lemma) *For every gauge $\delta(t)$ on a finite interval $[a, b]$ there is a δ -fine partition.*

Proof: Assume that $[a, b]$ lacks a δ -fine partition. Since we can construct a δ -fine partition of $[a, b]$ by appending a δ -fine partition of the half-interval $[(a+b)/2, b]$ to a δ -fine partition of the half-interval $[a, (a+b)/2]$, it follows that either $[a, (a+b)/2]$ or $[(a+b)/2, b]$ lacks a δ -fine partition. As in Example 2.3.1, we choose one of the half-intervals based on this failure and continue bisecting. This creates a nested sequence of intervals $[a_i, b_i]$ converging to a point x . If i is large enough, then $[a_i, b_i] \subset (x - \delta(x), x + \delta(x))$, and the interval $[a_i, b_i]$ with tag x is a δ -fine partition of itself. This contradicts the choice of $[a_i, b_i]$ and the assumption that the original interval $[a, b]$ lacks a δ -fine partition. ■

Before launching into our treatment of the gauge integral, we pause to gain some facility with gauge functions [108]. Here are three examples that illustrate their value.

Example 3.2.1 *A Gauge Proof of Weierstrass' Theorem*

Consider a real-valued continuous function $f(t)$ with domain $[a, b]$. Suppose that $f(t)$ does not attain its supremum on $[a, b]$. Then for each t there exists a point $x \in [a, b]$ with $f(t) < f(x)$. By continuity there exists $\delta(t) > 0$ such that $f(y) < f(x)$ for all $y \in [a, b]$ with $|y - t| < \delta(t)$. Using $\delta(t)$ as a gauge, select a δ -fine partition $a = s_0 < s_1 < \dots < s_n = b$ with tags $t_i \in [s_i, s_{i+1}]$ and designated points x_i satisfying $f(t_i) < f(x_i)$. Let x_{\max} be the point x_i having the largest value $f(x_i)$. Because x_{\max} lies in some interval $[s_i, s_{i+1}]$, we have $f(x_{\max}) < f(x_i)$. This contradiction discredits our assumption that $f(x)$ does not attain its supremum. A similar argument applies to the infimum. ■

Example 3.2.2 *A Gauge Proof of the Heine-Borel Theorem*

One can use Cousin's lemma to prove the Heine-Borel Theorem on the real line [278]. This theorem states that if C is a compact set contained in the union $\cup_{\alpha} O_{\alpha}$ of a collection of open sets O_{α} , then C is actually contained in the union of a finite number of the O_{α} . Suppose $C \subset [a, b]$. Define a gauge $\delta(t)$ so that the interval $(t - \delta(t), t + \delta(t))$ does not intersect C when $t \notin C$ and $(t - \delta(t), t + \delta(t))$ is contained in some O_{α} when $t \in C$. Based on $\delta(t)$, select a δ -fine partition $a = s_0 < s_1 < \dots < s_n = b$ with tags $t_i \in [s_i, s_{i+1}]$. By definition C is contained in the union $\cup_{t_i \in C} U_i$, where U_i is the set O_{α} covering t_i . The Heine-Borel theorem extends to compact sets in \mathbb{R}^n .

Example 3.2.3 *A Gauge Proof of the Intermediate Value Theorem*

Under the assumption of the previous example, let c be a number strictly between $f(a)$ and $f(b)$. If we assume that there is no $t \in [a, b]$ with $f(t) = c$, then there exists a positive number $\delta(t)$ such that either $f(x) < c$ for all

$x \in [a, b]$ with $|x - t| < \delta(t)$ or $f(x) > c$ for all $x \in [a, b]$ with $|x - t| < \delta(t)$. We now select a δ -fine partition $a = s_0 < s_1 < \cdots < s_n = b$ and observe that throughout each interval $[s_i, s_{i+1}]$ either $f(t) < c$ or $f(t) > c$. If to start $f(s_0) = f(a) < c$, then $f(s_1) < c$, which implies $f(s_2) < c$ and so forth until we get to $f(s_n) = f(b) < c$. This contradicts the assumption that c lies strictly between $f(a)$ and $f(b)$. With minor differences, the same proof works when $f(a) > c$. ■

In preparation for our next example and for the fundamental theorem of calculus later in this chapter, we must define derivatives. A real-valued function $f(t)$ defined on an interval $[a, b]$ possesses a derivative $f'(c)$ at $c \in [a, b]$ provided the limit

$$\lim_{t \rightarrow c} \frac{f(t) - f(c)}{t - c} = f'(c) \quad (3.1)$$

exists. At the endpoints a and b , the limit is necessarily one sided. Taking a sequential view of convergence, definition (3.1) means that for every sequence t_m converging to c we must have

$$\lim_{m \rightarrow \infty} \frac{f(t_m) - f(c)}{t_m - c} = f'(c).$$

In calculus, we learn the following rules for computing derivatives:

Proposition 3.2.2 *If $f(t)$ and $g(t)$ are differentiable functions on (a, b) , then*

$$\begin{aligned} [\alpha f(t) + \beta g(t)]' &= \alpha f'(t) + \beta g'(t) \\ [f(t)g(t)]' &= f'(t)g(t) + f(t)g'(t) \\ \left[\frac{1}{f(t)}\right]' &= -\frac{f'(t)}{f(t)^2}. \end{aligned}$$

In the third formula we must assume $f(t) \neq 0$. Finally, if $g(t)$ maps into the domain of $f(t)$, then the functional composition $f \circ g(t)$ has derivative

$$[f \circ g(t)]' = f' \circ g(t)g'(t).$$

Proof: We will prove the above sum, product, quotient, and chain rules in a broader context in Chap. 4. Our proofs will not rely on integration. ■

Example 3.2.4 *Strictly Increasing Functions*

Let $f(t)$ be a differentiable function on $[c, d]$ with strictly positive derivative. We now show that $f(t)$ is strictly increasing. For each $t \in [c, d]$ there exists $\delta(t) > 0$ such that

$$\frac{f(x) - f(t)}{x - t} > 0 \quad (3.2)$$

for all $x \in [a, b]$ with $|x - t| < \delta(t)$. According to Proposition 3.2.1, for any two points $a < b$ from $[c, d]$, there exists a δ -fine partition

$$a = s_0 < s_1 < \cdots < s_n = b$$

of $[a, b]$ with tags $t_i \in [s_i, s_{i+1}]$. In view of inequality (3.2), at least one of the two inequalities $f(s_i) \leq f(t_i) \leq f(s_{i+1})$ must be strict. Thus, the telescoping sum

$$f(b) - f(a) = \sum_{i=0}^{n-1} [f(s_{i+1}) - f(s_i)]$$

must be positive. ■

3.3 Definition and Basic Properties of the Integral

With later applications in mind, it will be convenient to define the gauge integral for vector-valued functions $f(x) : [a, b] \mapsto \mathbb{R}^n$. In this context, $f(x)$ is said to have integral \mathbf{I} if for every $\epsilon > 0$ there exists a gauge $\delta(x)$ on $[a, b]$ such that

$$\|S(f, \pi) - \mathbf{I}\| < \epsilon \tag{3.3}$$

for all δ -fine partitions π . Our first order of business is to check that the integral is unique whenever it exists. Thus, suppose that the vector \mathbf{J} is a second possible value of the integral. Given $\epsilon > 0$ choose gauges $\delta_{\mathbf{I}}(x)$ and $\delta_{\mathbf{J}}(x)$ leading to inequality (3.3). The minimum $\delta(x) = \min\{\delta_{\mathbf{I}}(x), \delta_{\mathbf{J}}(x)\}$ is also a gauge, and any partition π that is δ -fine is also $\delta_{\mathbf{I}}$ and $\delta_{\mathbf{J}}$ -fine. Hence,

$$\|\mathbf{I} - \mathbf{J}\| \leq \|\mathbf{I} - S(f, \pi)\| + \|S(f, \pi) - \mathbf{J}\| < 2\epsilon.$$

Since ϵ is arbitrary, $\mathbf{J} = \mathbf{I}$.

One can also define $f(x)$ to be integrable if its Riemann sums are Cauchy in an appropriate sense.

Proposition 3.3.1 (Cauchy criterion) *A function $f(x) : [a, b] \mapsto \mathbb{R}^n$ is integrable if and only if for every $\epsilon > 0$ there exists a gauge $\delta(x) > 0$ such that*

$$\|S(f, \pi_1) - S(f, \pi_2)\| < \epsilon \tag{3.4}$$

for any two δ -fine partitions π_1 and π_2 .

Proof: It is obvious that the Cauchy criterion is necessary for integrability. To show that it is sufficient, consider the sequence $\epsilon_m = m^{-1}$ and compatible sequence of gauges $\delta_m(x)$ determined by condition (3.4). We can force

the constraints $\delta_m(x) \leq \delta_{m-1}(x)$ to hold by inductively replacing $\delta_m(x)$ by $\min\{\delta_{m-1}(x), \delta_m(x)\}$ whenever needed. Now select a δ_m -fine partition π_m for each m . Because the gauge sequence $\delta_m(x)$ is decreasing, every partition π that is δ_m -fine is also δ_{m-1} -fine. Hence, the sequence of Riemann sums $S(f, \pi_m)$ is Cauchy and has a limit \mathbf{I} satisfying $\|S(f, \pi_m) - \mathbf{I}\| \leq m^{-1}$. Finally, given the potential integral \mathbf{I} , we take an arbitrary $\epsilon > 0$ and choose m so that $m^{-1} < \epsilon$. If π is δ_m -fine, then the inequality

$$\|S(f, \pi) - \mathbf{I}\| \leq \|S(f, \pi) - S(f, \pi_m)\| + \|S(f, \pi_m) - \mathbf{I}\| < 2\epsilon.$$

completes the proof. ■

For two integrable functions $f(x)$ and $g(x)$, the gauge integral inherits the linearity property

$$\int_a^b [\alpha f(x) + \beta g(x)] dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx$$

from its approximating Riemann sums. To prove this fact, take $\epsilon > 0$ and choose gauges $\delta_f(x)$ and $\delta_g(x)$ so that

$$\left\| S(f, \pi_f) - \int_a^b f(x) dx \right\| < \epsilon, \quad \left\| S(g, \pi_g) - \int_a^b g(x) dx \right\| < \epsilon$$

whenever π_f is δ_f -fine and π_g is δ_g -fine. If the tagged partition π is δ -fine for the gauge $\delta(x) = \min\{\delta_f(x), \delta_g(x)\}$, then

$$\begin{aligned} & \left\| S(\alpha f + \beta g, \pi) - \alpha \int_a^b f(x) dx - \beta \int_a^b g(x) dx \right\| \\ & \leq |\alpha| \left\| S(f, \pi) - \int_a^b f(x) dx \right\| + |\beta| \left\| S(g, \pi) - \int_a^b g(x) dx \right\| \\ & \leq (|\alpha| + |\beta|)\epsilon. \end{aligned}$$

The gauge integral also inherits obvious order properties. For example, $\int_a^b f(x) dx \geq 0$ whenever the integrand $f(x) \geq 0$ for all $x \in [a, b]$. In this case, the inequality $|S(f, \pi) - \int_a^b f(x) dx| < \epsilon$ implies

$$0 \leq S(f, \pi) \leq \int_a^b f(x) dx + \epsilon.$$

Since ϵ can be made arbitrarily small for $f(x)$ integrable, it follows that $\int_a^b f(x) dx \geq 0$. This nonnegativity property translates into the order property

$$\int_a^b f(x) dx \leq \int_a^b g(x) dx$$

for two integrable functions $f(x) \leq g(x)$. In particular, when both $f(x)$ and $|f(x)|$ are both integrable, we have

$$\left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx.$$

For vector-valued functions, the analogous rule

$$\left\| \int_a^b f(x) dx \right\| \leq \int_a^b \|f(x)\| dx \quad (3.5)$$

is also inherited from the approximating Riemann sums. The reader can easily supply the proof using the triangle inequality of the Euclidean norm. It does not take much imagination to extend the definition of the gauge integral to matrix-valued functions, and inequality (3.5) applies in this setting as well.

One of the nicest features of the gauge integral is that one can perturb an integrable function at a countable number of points without changing the value of its integral. This property fails for the Riemann integral but is exhibited by the Lebesgue integral. To validate the property, it suffices to prove that a function that equals $\mathbf{0}$ except at a countable number of points has integral $\mathbf{0}$. Suppose $f(x)$ is such a function with exceptional points x_1, x_2, \dots and corresponding exceptional values $\mathbf{f}_1, \mathbf{f}_2, \dots$. We now define a gauge $\delta(x)$ with value 1 on the nonexceptional points and values

$$\delta(x_j) = \frac{\epsilon}{2^{j+2}[\|\mathbf{f}_j\| + 1]}$$

at the exceptional points. If π is a δ -fine partition, then x_j can serve as a tag for at most two intervals $[s_i, s_{i+1}]$ of π and each such interval has length less than $2\delta(x_j)$. It follows that

$$\|S(f, \pi)\| \leq 2 \sum_j \|f(x_j)\| \frac{2\epsilon}{2^{j+2}[\|\mathbf{f}_j\| + 1]} \leq \epsilon \sum_{j=1}^{\infty} \frac{1}{2^j} = \epsilon$$

and therefore that $\int_a^b f(x) dx = \mathbf{0}$.

In practice, the interval additivity rule

$$\int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx \quad (3.6)$$

is obviously desirable. There are three separate issues in proving it. First, given the existence of the integral over $[a, c]$, do the integrals over $[a, b]$ and $[b, c]$ exist? Second, if the integrals over $[a, b]$ and $[b, c]$ exist, does the integral over $[a, c]$ exist? Third, if the integrals over $[a, b]$ and $[b, c]$ exist, are they additive? The first question is best approached through Proposition 3.3.1. For $\epsilon > 0$ there exists a gauge $\delta(x)$ such that

$$\|S(f, \pi_1) - S(f, \pi_2)\| < \epsilon$$

for any two δ -fine partitions π_1 and π_2 of $[a, c]$. Given $\delta(x)$, take any two δ -fine partitions γ_1 and γ_2 of $[a, b]$ and a single δ -fine partition ω of $[b, c]$. The concatenated partitions $\gamma_1 \cup \omega$ and $\gamma_2 \cup \omega$ are δ -fine throughout $[a, c]$ and satisfy

$$\|S(f, \gamma_1) - S(f, \gamma_2)\| = \|S(f, \gamma_1 \cup \omega) - S(f, \gamma_2 \cup \omega)\| < \epsilon.$$

According to the Cauchy criterion, the integral over $[a, b]$ therefore exists. A similar argument implies that the integral over $[b, c]$ also exists. Finally, the combination of these results shows that the integral exists over any interval $[u, v]$ contained within $[a, b]$.

For the converse, choose gauges $\delta_1(x)$ on $[a, b]$ and $\delta_2(x)$ on $[b, c]$ so that

$$\left\| S(f, \gamma) - \int_a^b f(x) dx \right\| < \epsilon, \quad \left\| S(f, \omega) - \int_b^c f(x) dx \right\| < \epsilon$$

for any δ_1 -fine partition γ of $[a, b]$ and any δ_2 -fine partition ω of $[b, c]$. The concatenated partition $\pi = \gamma \cup \omega$ satisfies

$$\begin{aligned} & \left\| S(f, \pi) - \int_a^b f(x) dx - \int_b^c f(x) dx \right\| \\ & \leq \left\| S(f, \gamma) - \int_a^b f(x) dx \right\| + \left\| S(f, \omega) - \int_b^c f(x) dx \right\| \\ & < 2\epsilon \end{aligned} \quad (3.7)$$

because the Riemann sums satisfy $S(f, \pi) = S(f, \gamma) + S(f, \omega)$. This suggests defining a gauge $\delta(x)$ equal to $\delta_1(x)$ on $[a, b]$ and equal to $\delta_2(x)$ on $[b, c]$. The problem with this tactic is that some partitions of $[a, c]$ do not split at b . However, we can ensure a split by redefining $\delta(x)$ by

$$\tilde{\delta}(x) = \begin{cases} \min\{\delta_1(b), \delta_2(b)\} & x = b \\ \min\{\delta(x), \frac{1}{2}|x - b|\} & x \neq b. \end{cases}$$

This forces b to be the tag of its assigned interval, and we can if needed split this interval at b and retain b as tag of both subintervals. With $\delta(x)$ amended in this fashion, any δ -fine partition π can be viewed as a concatenated partition $\gamma \cup \omega$ splitting at b . As such π obeys inequality (3.7). This argument simultaneously proves that the integral over $[a, c]$ exists and satisfies the additivity property (3.6)

If the function $f(x)$ is vector-valued with n components, then the integrability of $f(x)$ should imply the integrability of each its components $f_i(x)$. Furthermore, we should be able to write

$$\int_a^b f(x) dx = \begin{pmatrix} \int_a^b f_1(x) dx \\ \vdots \\ \int_a^b f_n(x) dx \end{pmatrix}.$$

Conversely, if its components are integrable, then $f(x)$ should be integrable as well. The inequalities

$$\|S(f, \pi) - \mathbf{I}\| \leq \sum_{i=1}^n |S(f_i, \pi) - I_i| \leq \sqrt{n} \|S(f, \pi) - \mathbf{I}\|.$$

based on Example 2.5.6 and Problem 3 of Chap. 2 are instrumental in proving this logical equivalence. Given that we can integrate component by component, for the remainder of this chapter we will deal exclusively with real-valued functions.

We have not actually shown that any function is integrable. The most obvious possibility is a constant. Fortunately, it is trivial to demonstrate that

$$\int_a^b c \, dx = c(b - a).$$

Step functions are one rung up the hierarchy of functions. If

$$f(x) = \sum_{i=0}^{n-1} c_i 1_{(s_i, s_{i+1}]}(x)$$

for $a = s_0 < s_1 < \dots < s_n = b$, then our nascent theory allows us to evaluate

$$\int_a^b f(x) \, dx = \sum_{i=0}^{n-1} \int_{s_i}^{s_{i+1}} c_i \, dx = \sum_{i=0}^{n-1} c_i (s_{i+1} - s_i).$$

This fact and the next technical proposition turn out to be the key to showing that continuous functions are integrable.

Proposition 3.3.2 *Let $f(x)$ be a function with domain $[a, b]$. Suppose for every $\epsilon > 0$ there exist two integrable functions $g(x)$ and $h(x)$ satisfying $g(x) \leq f(x) \leq h(x)$ for all x and*

$$\int_a^b h(x) \, dx \leq \int_a^b g(x) \, dx + \epsilon.$$

Then $f(x)$ is integrable.

Proof: For $\epsilon > 0$, choose gauges $\delta_g(x)$ and $\delta_h(x)$ on $[a, b]$ so that

$$\left| S(g, \pi_g) - \int_a^b g(x) \, dx \right| < \epsilon, \quad \left| S(h, \pi_h) - \int_a^b h(x) \, dx \right| < \epsilon$$

for any δ_g -fine partition π_g and any δ_h -fine partition π_h . If π is a δ -fine partition for $\delta(x) = \min\{\delta_g(x), \delta_h(x)\}$, then the inequalities

$$\begin{aligned} \int_a^b g(x) dx - \epsilon &< S(g, \pi) \\ &\leq S(f, \pi) \\ &\leq S(h, \pi) \\ &< \int_a^b h(x) dx + \epsilon \\ &\leq \int_a^b g(x) dx + 2\epsilon \end{aligned}$$

trap $S(f, \pi)$ in an interval of length 3ϵ . Because the Riemann sum $S(f, \gamma)$ for any other δ -fine partition γ is trapped in the same interval, the integral of $f(x)$ exists by the Cauchy criterion. ■

Proposition 3.3.3 *Every continuous function $f(x)$ on $[a, b]$ is integrable.*

Proof: In view of the uniform continuity of $f(x)$ on $[a, b]$, for every $\epsilon > 0$ there exists a $\delta > 0$ with $|f(x) - f(y)| < \epsilon$ when $|x - y| < \delta$. For the constant gauge $\delta(x) = \delta$ and a corresponding δ -fine partition π with mesh points s_0, \dots, s_n , let m_i be the minimum and M_i be the maximum of $f(x)$ on $[s_i, s_{i+1}]$. The step functions

$$g(x) = \sum_{i=1}^n m_i 1_{(s_i, s_{i+1}]}(x), \quad h(x) = \sum_{i=1}^n M_i 1_{(s_i, s_{i+1}]}(x)$$

then satisfy $g(x) \leq f(x) \leq h(x)$ except at the single point a . Furthermore,

$$\begin{aligned} \int_a^b h(x) dx - \int_a^b g(x) dx &\leq \epsilon \sum_{i=1}^n (s_{i+1} - s_i) \\ &= \epsilon(b - a). \end{aligned}$$

Application of Proposition 3.3.2 now completes the proof. ■

3.4 The Fundamental Theorem of Calculus

The fundamental theorem of calculus divides naturally into two parts. For the gauge integral, the first and more difficult part is easily proved by invoking what is called the straddle inequality. Let $f(x)$ be differentiable at the point $t \in [a, b]$. Then there exists $\delta(t) > 0$ such that

$$\left| \frac{f(x) - f(t)}{x - t} - f'(t) \right| < \epsilon$$

for all $x \in [a, b]$ with $|x - t| < \delta(t)$. If $u < t < v$ are two points straddling t and located in $[a, b] \cap (t - \delta(t), t + \delta(t))$, then

$$\begin{aligned} |f(v) - f(u) - f'(t)(v - u)| &\leq |f(v) - f(t) - f'(t)(v - t)| \\ &\quad + |f(t) - f(u) - f'(t)(t - u)| \\ &\leq \epsilon(v - t) + \epsilon(t - u) \\ &= \epsilon(v - u). \end{aligned} \tag{3.8}$$

Inequality (3.8) also clearly holds when either $u = t$ or $v = t$.

Proposition 3.4.1 (Fundamental Theorem I) *If $f(x)$ is differentiable throughout $[a, b]$, then*

$$\int_a^b f'(x) dx = f(b) - f(a).$$

Proof: Using the gauge $\delta(t)$ figuring in the straddle inequality (3.8), select a δ -fine partition π with mesh points $a = s_0 < s_1 < \dots < s_n = b$ and tags $t_i \in [s_i, s_{i+1}]$. Application of the inequality and telescoping yield

$$\begin{aligned} |f(b) - f(a) - S(f', \pi)| &= \left| \sum_{i=0}^{n-1} [f(s_{i+1}) - f(s_i) - f'(t_i)(s_{i+1} - s_i)] \right| \\ &\leq \sum_{i=0}^{n-1} |f(s_{i+1}) - f(s_i) - f'(t_i)(s_{i+1} - s_i)| \\ &\leq \sum_{i=0}^{n-1} \epsilon(s_{i+1} - s_i) \\ &= \epsilon(b - a). \end{aligned}$$

This demonstrates that $f'(x)$ has integral $f(b) - f(a)$. ■

The first half of the fundamental theorem remains valid for a continuous function $f(x)$ that is differentiable except on a countable set N [250]. Since changing an integrand at a countable number of points does not alter its integral, it suffices to prove that

$$f(b) - f(a) = \int_a^b g(t) dt, \quad \text{where } g(t) = \begin{cases} 0 & t \in N \\ f'(t) & t \notin N. \end{cases}$$

Suppose $\epsilon > 0$ is given. For $t \notin N$ define the gauge value $\delta(t)$ to satisfy the straddle inequality. Enumerate the points t_j of N , and define $\delta(t_j) > 0$ so that $|f(t_j) - f(t_j + s)| < 2^{-j-2}\epsilon$ whenever $|s| < \delta(t_j)$. Now select a δ -fine partition π with mesh points $a = s_0 < s_1 < \dots < s_n = b$ and tags $r_i \in [s_i, s_{i+1}]$. Break the sum

$$f(b) - f(a) - S(g, \pi) = \sum_{i=0}^{n-1} [f(s_{i+1}) - f(s_i) - g(r_i)(s_{i+1} - s_i)]$$

into two parts. Let S' denote the sum of the terms with tags $r_i \notin N$, and let S'' denote the sum of the terms with tags $r_i \in N$. As noted earlier, $|S'| \leq \epsilon(b-a)$. Because a tag is attached to at most two subintervals, the second sum satisfies

$$\begin{aligned} |S''| &\leq \sum_{r_i \in N} |f(s_{i+1}) - f(s_i)| \\ &\leq \sum_{r_i \in N} [|f(s_{i+1}) - f(r_i)| + |f(r_i) - f(s_i)|] \\ &\leq 2 \sum_{j=1}^{\infty} 22^{-j-2} \epsilon = \epsilon. \end{aligned}$$

It follows that $|S' + S''| \leq \epsilon(b-a+1)$ and therefore that the stated integral exists and equals $f(b) - f(a)$.

In demonstrating the second half of the fundamental theorem, we will implicitly use the standard convention

$$\int_d^c f(x) dx = - \int_c^d f(x) dx$$

for $c < d$. This convention will also be in force in proving the substitution formula.

Proposition 3.4.2 (Fundamental Theorem II) *If a function $f(x)$ is integrable on $[a, b]$, then its indefinite integral*

$$F(t) = \int_a^t f(x) dx$$

has derivative $F'(t) = f(t)$ at any point t where $f(x)$ is continuous. The derivative is taken as one sided if $t = a$ or $t = b$.

Proof: In deriving the interval additivity rule (3.6), we showed that the integral $F(t)$ exists. At a point t where $f(x)$ is continuous, for any $\epsilon > 0$ there is a $\delta > 0$ such that $-\epsilon < f(x) - f(t) < \epsilon$ when $|x - t| < \delta$ and $x \in [a, b]$. Hence, the difference

$$\frac{F(t+s) - F(t)}{s} - f(t) = \frac{1}{s} \int_t^{t+s} [f(x) - f(t)] dx$$

is less than ϵ and greater than $-\epsilon$ for $|s| < \delta$. In the limit as s tends to 0, we recover $F'(t) = f(t)$. ■

The fundamental theorem of calculus has several important corollaries. These are covered in the next three propositions on the substitution rule, integration by parts, and finite Taylor expansions.

Proposition 3.4.3 (Substitution Rule) *Suppose $f(x)$ is differentiable on $[a, b]$, $g(x)$ is differentiable on $[c, d]$, and the image of $[c, d]$ under $g(x)$ is contained within $[a, b]$. Then*

$$\int_{g(c)}^{g(d)} f'(y) dy = \int_c^d f'[g(x)]g'(x) dx.$$

Proof: Part I of the fundamental theorem and the chain rule identity

$$\{f[g(x)]\}' = f'[g(x)]g'(x)$$

imply that both integrals have value $f[g(d)] - f[g(c)]$. ■

Proposition 3.4.4 (Integration by Parts) *Suppose $f(x)$ and $g(x)$ are differentiable on $[a, b]$. Then $f'(x)g(x)$ is integrable on $[a, b]$ if and only if $f(x)g'(x)$ is integrable on $[a, b]$. Furthermore, the two integrals are related by the identity*

$$\int_a^b f'(x)g(x) dx + \int_a^b f(x)g'(x) dx = f(b)g(b) - f(a)g(a),$$

Proof: The product rule for derivatives is

$$[f(x)g(x)]' = f'(x)g(x) + f(x)g'(x).$$

If two of three members of this identity are integrable, then the third is as well. Since part I of the fundamental theorem entails

$$\int_a^b [f(x)g(x)]' dx = f(b)g(b) - f(a)g(a),$$

the proposition follows. ■

The derivative of a function may itself be differentiable. Indeed, it makes sense to speak of the k th-order derivative of a function $f(x)$ if $f(x)$ is sufficiently smooth. Traditionally, the second-order derivative is denoted $f''(x)$ and an arbitrary k th-order derivative by $f^{(k)}(x)$. We can use these extra derivatives to good effect in approximating $f(x)$ locally. The next proposition makes this clear and offers an explicit estimate of the error in a finite Taylor expansion of $f(x)$.

Proposition 3.4.5 (Taylor Expansion) *Suppose $f(x)$ has a derivative of order $k+1$ on an open interval around the point y . Then for all x in the interval, we have*

$$f(x) = f(y) + \sum_{j=1}^k \frac{1}{j!} f^{(j)}(y)(x-y)^j + R_k(x), \quad (3.9)$$

where the remainder

$$R_k(x) = \frac{(x-y)^{k+1}}{k!} \int_0^1 f^{(k+1)}[y+t(x-y)](1-t)^k dt.$$

If $|f^{(k+1)}(z)| \leq b$ for all z between x and y , then

$$|R_k(x)| \leq \frac{b|x-y|^{k+1}}{(k+1)!}. \quad (3.10)$$

Proof: When $k = 0$, the Taylor expansion (3.9) reads

$$f(x) = f(y) + (x-y) \int_0^1 f'[y+t(x-y)]dt$$

and follows from the fundamental theorem of calculus and the chain rule. Induction and the integration-by-parts formula

$$\begin{aligned} & \int_0^1 f^{(k)}[y+t(x-y)](1-t)^{k-1} dt \\ &= -\frac{1}{k} f^{(k)}[y+t(x-y)](1-t)^k \Big|_0^1 \\ & \quad + \frac{x-y}{k} \int_0^1 f^{(k+1)}[y+t(x-y)](1-t)^k dt \\ &= \frac{1}{k} f^{(k)}(y) + \frac{x-y}{k} \int_0^1 f^{(k+1)}[y+t(x-y)](1-t)^k dt \end{aligned}$$

now validate the general expansion (3.9). The error estimate follows directly from the bound $|f^{(k+1)}(z)| \leq b$ and the integral

$$\int_0^1 (1-t)^k dt = \frac{1}{k+1}.$$

3.5 More Advanced Topics in Integration

Within the confines of a single chapter, it is impossible to develop rigorously all of the properties of the gauge integral. In this section we will discuss briefly four topics: (a) integrals over unbounded intervals, (b) improper integrals and Hake's theorem, (c) the interchange of limits and integrals, and (d) multidimensional integrals and Fubini's theorem.

Defining the integral of a function over an unbounded interval requires several minor adjustments. First, the real line is extended to include the points $\pm\infty$. Second, a gauge function $\delta(x)$ is now viewed as mapping x to an open interval containing x . The associated interval may be infinite; indeed, it must be infinite if x equals $\pm\infty$. In a δ -fine partition π , the

interval I_j containing the tag x_j is contained in $\delta(x_j)$. The length of an infinite interval I_j is defined to be 0 in an approximating Riemann sum $S(f, \pi)$ to avoid infinite contributions to the sum. Likewise, the integrand $f(x)$ is assigned the value 0 at $x = \pm\infty$.

This extended definition carries with it all the properties we expect. Its most remarkable consequence is that it obliterates the distinction between proper and improper integrals. Hake's theorem provides the link. If we allow a and b to be infinite as well as finite, then Hake's theorem says a function $f(x)$ is integrable over (a, b) if and only if either of the two limits

$$\lim_{c \rightarrow a} \int_c^b f(x) dx \quad \text{or} \quad \lim_{c \rightarrow b} \int_a^c f(x) dx$$

exists. If either limit exists, then $\int_a^b f(x) dx$ equals that limit. For instance, the integral

$$\int_1^\infty \frac{1}{x^2} dx = \lim_{c \rightarrow \infty} \int_1^c \frac{1}{x^2} dx = \lim_{c \rightarrow \infty} -\frac{1}{x} \Big|_1^c = 1$$

exists and has the indicated limit by this reasoning.

Example 3.5.1 *Existence of $\int_0^\infty \text{sinc}(x) dx$*

Consider the integral of $\text{sinc}(x) = \sin(x)/x$ over the interval $(0, \infty)$. Because $\text{sinc}(x)$ is continuous throughout $[0, 1]$ with limit 1 as x approaches 0, the integral over $[0, 1]$ is defined. Hake's theorem and integration by parts show that the integral

$$\begin{aligned} \int_1^\infty \frac{\sin x}{x} dx &= \lim_{c \rightarrow \infty} \int_1^c \frac{\sin x}{x} dx \\ &= \lim_{c \rightarrow \infty} \left(-\frac{\cos x}{x} \Big|_1^c - \int_1^c \frac{\cos x}{x^2} dx \right) \\ &= \cos 1 - \int_1^\infty \frac{\cos x}{x^2} dx \end{aligned}$$

exists provided the integral of $x^{-2} \cos x$ exists over $(1, \infty)$. We will demonstrate this fact in a moment. If we accept it, then it is clear that the integral of $\text{sinc}(x)$ over $(0, \infty)$ exists as well. As we shall find in Example 3.5.4, this integral equals $\pi/2$. In contrast to these positive results, $\text{sinc}(x)$ is not absolutely integrable over $(0, \infty)$. Finally, we note in passing that the substitution rule gives

$$\int_0^\infty \frac{\sin cx}{x} dx = \int_0^\infty \frac{\sin y}{c^{-1}y} c^{-1} dy = \int_0^\infty \frac{\sin y}{y} dy = \frac{\pi}{2}.$$

for any $c > 0$. ■

We now ask under what circumstances the formula

$$\lim_{n \rightarrow \infty} \int_a^b f_n(x) dx = \int_a^b \lim_{n \rightarrow \infty} f_n(x) dx \quad (3.11)$$

is valid. The two relevant theorems permitting the interchange of limits and integrals are the monotone convergence theorem and the dominated convergence theorem. In the monotone convergence theorem, we are given an increasing sequence $f_n(x)$ of integrable functions that converge to a finite limit for each x . Formula (3.11) is true in this setting provided

$$\sup_n \int_a^b f_n(x) dx < \infty.$$

In the dominated convergence theorem, we assume the sequence $f_n(x)$ is trapped between two integrable functions $g(x)$ and $h(x)$ in the sense that

$$g(x) \leq f_n(x) \leq h(x)$$

for all n and x . If $\lim_{n \rightarrow \infty} f_n(x)$ exists in this setting, then the interchange (3.11) is allowed. The choices

$$f_n(x) = 1_{[1,n]}(x)x^{-2} \cos x, \quad g(x) = -x^{-2}, \quad h(x) = x^{-2}$$

in the dominated convergence theorem validate the existence of

$$\int_1^{\infty} x^{-2} \cos x dx = \lim_{n \rightarrow \infty} \int_1^n x^{-2} \cos x dx.$$

We now consider two more substantive applications of the monotone and dominated convergence theorems.

Example 3.5.2 *Johann Bernoulli's Integral*

As example of delicate maneuvers in integration, consider the integral

$$\begin{aligned} \int_0^1 \frac{1}{x^x} dx &= \int_0^1 e^{-x \ln x} dx \\ &= \int_0^1 \sum_{n=0}^{\infty} \frac{(-x \ln x)^n}{n!} dx \\ &= \sum_{n=0}^{\infty} \frac{1}{n!} \int_0^1 (-x \ln x)^n dx. \end{aligned}$$

The reader will notice the application of the monotone convergence theorem in passing from the second to the third line above. Further progress can be made by applying the integration by parts result

$$\begin{aligned} \int_0^1 x^m \ln^n x dx &= -\frac{n}{m+1} \int_0^1 x^{m+1} \frac{\ln^{n-1} x}{x} dx \\ &= -\frac{n}{m+1} \int_0^1 x^m \ln^{n-1} x dx \end{aligned}$$

recursively to evaluate

$$\int_0^1 (-x \ln x)^n dx = \frac{n!}{(n+1)^n} \int_0^1 x^n dx = \frac{n!}{(n+1)^{n+1}}.$$

The pleasant surprise

$$\int_0^1 \frac{1}{x^x} dx = \sum_{n=0}^{\infty} \frac{1}{(n+1)^{n+1}}$$

emerges. ■

Example 3.5.3 *Competing Definitions of the Gamma Function*

The dominated convergence theorem allows us to derive Gauss's representation

$$\Gamma(z) = \lim_{n \rightarrow \infty} \frac{n! n^z}{z(z+1) \cdots (z+n)}$$

of the gamma function from Euler's representation

$$\Gamma(z) = \int_0^{\infty} x^{z-1} e^{-x} dx.$$

As students of statistics are apt to know from their exposure to the beta distribution, repeated integration by parts and the fundamental theorem of calculus show that

$$\int_0^1 x^{z-1} (1-x)^n dx = \frac{n!}{z(z+1) \cdots (z+n)}.$$

The substitution rule yields

$$n^z \int_0^1 x^{z-1} (1-x)^n dx = \int_0^n y^{z-1} \left(1 - \frac{y}{n}\right)^n dy.$$

Thus, it suffices to prove that

$$\int_0^{\infty} x^{z-1} e^{-x} dx = \lim_{n \rightarrow \infty} \int_0^n y^{z-1} \left(1 - \frac{y}{n}\right)^n dy.$$

Given the limit

$$\lim_{n \rightarrow \infty} \left(1 - \frac{y}{n}\right)^n = e^{-y},$$

we need an integrable function $h(y)$ that dominates the nonnegative sequence

$$f_n(y) = 1_{[0,n]}(y) y^{z-1} \left(1 - \frac{y}{n}\right)^n$$

from above in order to apply the dominated convergence theorem. In light of the inequality

$$\left(1 - \frac{y}{n}\right)^n \leq e^{-y},$$

the function $h(y) = y^{z-1}e^{-y}$ will serve. ■

Finally, the gauge integral extends to multiple dimensions, where a version of Fubini's theorem holds for evaluating multidimensional integrals via iterated integrals [278]. Consider a function $f(\mathbf{x}, \mathbf{y})$ defined over the Cartesian product $H \times K$ of two multidimensional intervals H and K . The intervals in question can be bounded or unbounded. If $f(\mathbf{x}, \mathbf{y})$ is integrable over $H \times K$, then Fubini's theorem asserts that the integrals $\int_H f(\mathbf{x}, \mathbf{y}) d\mathbf{x}$ and $\int_K f(\mathbf{x}, \mathbf{y}) d\mathbf{y}$ exist and can be integrated over the remaining variable to give the full integral. In symbols,

$$\int_{H \times K} f(\mathbf{x}, \mathbf{y}) d\mathbf{x} d\mathbf{y} = \int_K \left[\int_H f(\mathbf{x}, \mathbf{y}) d\mathbf{x} \right] d\mathbf{y} = \int_H \left[\int_K f(\mathbf{x}, \mathbf{y}) d\mathbf{y} \right] d\mathbf{x}.$$

Conversely, if either iterated integral exists, one would like to conclude that the full integral exists as well. This is true whenever $f(\mathbf{x}, \mathbf{y})$ is nonnegative. Unfortunately, it is false in general, and two additional hypotheses introduced by Tonelli are needed to rescue the situation. One hypothesis is that $f(\mathbf{x}, \mathbf{y})$ is measurable. Measurability is a technical condition that holds except for very pathological functions. The other hypothesis is that $|f(\mathbf{x}, \mathbf{y})| \leq g(\mathbf{x}, \mathbf{y})$ for some nonnegative function $g(\mathbf{x}, \mathbf{y})$ for which the iterated integral exists. This domination condition is shared with the dominated convergence theorem and forces $f(\mathbf{x}, \mathbf{y})$ to be absolutely integrable.

Example 3.5.4 *Evaluation of $\int_0^\infty \text{sinc}(x) dx$*

According to Fubini's theorem

$$\int_0^n \int_0^{n\pi} e^{-xy} \sin x dx dy = \int_0^{n\pi} \int_0^n e^{-xy} \sin x dy dx. \quad (3.12)$$

The second of these iterated integrals

$$\int_0^{n\pi} \int_0^n e^{-xy} \sin x dy dx = \int_0^{n\pi} \frac{\sin x}{x} dx - \int_0^{n\pi} e^{-nx} \frac{\sin x}{x} dx$$

tends to $\int_0^\infty \text{sinc}(x) dx$ as n tends to ∞ by a combination of Hake's theorem and the dominated convergence theorem. The inner integral of the left iterated integral in (3.12) equals

$$\begin{aligned} \int_0^{n\pi} e^{-xy} \sin x dx &= -e^{-xy} \cos x \Big|_0^{n\pi} - ye^{-xy} \sin x \Big|_0^{n\pi} \\ &\quad - y^2 \int_0^{n\pi} e^{-xy} \sin x dx \\ &= 1 - e^{-n\pi y} \cos n\pi - y^2 \int_0^{n\pi} e^{-xy} \sin x dx \end{aligned}$$

after two integrations by parts. It follows that

$$\int_0^{n\pi} e^{-xy} \sin x \, dx = \frac{1 - e^{-n\pi y} \cos n\pi}{1 + y^2}.$$

Finally, application of the dominated convergence theorem gives

$$\begin{aligned} \lim_{n \rightarrow \infty} \int_0^n \frac{1 - e^{-n\pi y} \cos n\pi}{1 + y^2} \, dy &= \int_0^\infty \frac{1}{1 + y^2} \, dy \\ &= \frac{\pi}{2}. \end{aligned}$$

Equating the limits of the right and left hand sides of the identity (3.12) therefore [278] yields the value of $\pi/2$ for $\int_0^\infty \operatorname{sinc}(x) \, dx$. ■

3.6 Problems

1. Give an alternative proof of Cousin's lemma by letting y be the supremum of the set of $x \in [a, b]$ such that $[a, x]$ possesses a δ -fine partition.
2. Use Cousin's lemma to prove that a continuous function $f(x)$ defined on an interval $[a, b]$ is uniformly continuous there [108]. (Hint: Given $\epsilon > 0$ define a gauge $\delta(x)$ by the requirement that $|f(y) - f(x)| < \frac{1}{2}\epsilon$ for all $y \in [a, b]$ with $|y - x| < 2\delta(x)$.)
3. A possibly discontinuous function $f(x)$ has one-sided limits at each point $x \in [a, b]$. Show by Cousin's lemma that $f(x)$ is bounded on $[a, b]$.
4. Suppose $f(x)$ has a nonnegative derivative $f'(x)$ throughout $[a, b]$. Prove that $f(x)$ is nondecreasing on $[a, b]$. Also prove that $f(x)$ is constant on $[a, b]$ if and only if $f'(x) = 0$ for all x . (Hint: These yield easily to the fundamental theorem of calculus. Alternatively for the first assertion, consider the function

$$f_\epsilon(x) = f(x) + \epsilon x$$

for $\epsilon > 0$.)

5. Using only the definition of the gauge integral, demonstrate that

$$\int_a^b f(t) \, dt = \int_{-b}^{-a} f(-t) \, dt$$

when either integral exists.

6. Based on the standard definition of the natural logarithm

$$\ln y = \int_1^y \frac{1}{x} dx,$$

prove that $\ln yz = \ln y + \ln z$ for all positive arguments y and z . Use this property to verify that $\ln y^{-1} = -\ln y$ and that $\ln y^r = r \ln y$ for every rational number r .

7. Apply Proposition 3.3.2 and demonstrate that every monotonic function defined on an interval $[a, b]$ is integrable on that interval.
8. Let $f(x)$ be a continuous real-valued function on $[a, b]$. Show that there exists $c \in [a, b]$ with

$$\int_a^b f(x) dx = f(c)(b - a).$$

9. In the Taylor expansion of Proposition 3.4.5, suppose $f^{(k+1)}(x)$ is continuous. Show that we can replace the remainder by

$$R_k(x) = \frac{(x - y)^{k+1}}{(k + 1)!} f^{(k+1)}(z)$$

for some z between x and y .

10. Suppose that $f(x)$ is infinitely differentiable and that c and r are positive numbers. If $|f^{(k)}(x)| \leq ck!r^k$ for all x near y and all nonnegative integers k , then use Proposition 3.4.5 to show that

$$f(x) = \sum_{k=0}^{\infty} \frac{f^{(k)}(y)}{k!} (x - y)^k$$

near y . Explicitly determine the infinite Taylor series expansion of the function $f(x) = (1 + x)^{-1}$ around $x = 0$ and justify its convergence.

11. Suppose the nonnegative continuous function $f(x)$ satisfies

$$\int_a^b f(x) dx = 0.$$

Prove that $f(x)$ is identically 0 on $[a, b]$.

12. Consider the function

$$f(x) = \begin{cases} x^2 \sin(x^{-2}) & x \neq 0 \\ 0 & x = 0. \end{cases}$$

Show that $\int_0^1 f'(x) dx = \sin(1)$ and $\lim_{t \downarrow 0} \int_t^1 |f'(x)| dx = \infty$. Hence, $f'(x)$ is integrable but not absolutely integrable on $[0, 1]$.

13. Prove that

$$\int_0^{\infty} x^{\alpha} e^{-x^{\beta}} dx = \frac{1}{\beta} \Gamma\left(\frac{\alpha+1}{\beta}\right)$$

for α and β positive [82].

14. Justify the formula

$$\int_0^1 \frac{\ln(1-x)}{x} dx = -\sum_{n=1}^{\infty} \frac{1}{n^2}.$$

15. Show that

$$\int_0^{\infty} \frac{x^{z-1}}{e^x - 1} dx = \zeta(z)\Gamma(z),$$

where $\zeta(z) = \sum_{n=1}^{\infty} n^{-z}$.

16. Prove that the functions

$$\begin{aligned} f(x) &= \int_1^{\infty} \frac{\sin t}{x^2 + t^2} dt \\ g(x) &= \int_0^{\infty} e^{-xt} \cos t dt, \quad x > 0, \end{aligned}$$

are continuous.

17. Let $f_n(x)$ be a sequence of integrable functions on $[a, b]$ that converges uniformly to $f(x)$. Demonstrate that $f(x)$ is integrable and satisfies

$$\lim_{n \rightarrow \infty} \int_a^b f_n(x) dx = \int_a^b f(x) dx.$$

(Hints: For $\epsilon > 0$ small take n large enough so that

$$f_n(x) - \frac{\epsilon}{2(b-a)} \leq f(x) \leq f_n(x) + \frac{\epsilon}{2(b-a)}$$

for all x .)

18. Let p and q be positive integers. Justify the series expansion

$$\int_0^1 \frac{x^{p-1}}{1+x^q} dx = \sum_{n=0}^{\infty} \frac{(-1)^n}{p+nq}$$

by the monotone convergence theorem. Be careful since the series does not converge absolutely [278].

19. Suppose $f(x)$ is a continuous function on \mathbb{R} . Demonstrate that the sequence

$$f_n(x) = \frac{1}{n} \sum_{k=0}^{n-1} f\left(x + \frac{k}{n}\right)$$

converges uniformly to a continuous function on every finite interval $[a, b]$ [69].

20. Prove that

$$\int_0^1 \frac{x^b - x^a}{\ln x} dx = \ln \frac{b+1}{a+1}$$

for $0 < a < b$ [278] by showing that both sides equal the double integral

$$\int_{[0,1] \times [a,b]} x^y dx dy.$$

21. Integrate the function

$$f(x, y) = \frac{y^2 - x^2}{(x^2 + y^2)^2}$$

over the unit square $[0, 1] \times [0, 1]$. Show that the two iterated integrals disagree, and explain why Fubini's theorem fails.

22. Suppose the two partial derivatives $\frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x})$ and $\frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x})$ exist and are continuous in a neighborhood of a point $\mathbf{y} \in \mathbb{R}^2$. Show that they are equal at the point. (Hints: If they are not equal, take a small box around the point where their difference has constant sign. Now apply Fubini's theorem.)

23. Demonstrate that

$$\int_0^\infty e^{-x^2} dx = \frac{\sqrt{\pi}}{2}$$

by evaluating the integral of $f(\mathbf{y}) = y_2 e^{-(1+y_1^2)y_2^2}$ over the rectangle $(0, \infty) \times (0, \infty)$.

4

Differentiation

4.1 Introduction

Differentiation and integration are the two pillars on which all of calculus rests. For real-valued functions of a real variable, all of the major issues surrounding differentiation were settled long ago. For multivariate differentiation, there are still some subtleties and snares. We adopt a definition of differentiability that avoids most of the pitfalls and makes differentiation of vectors and matrices relatively painless. In later chapters, this definition also improves the clarity of exposition.

The main theme of differentiation is the short-range approximation of curved functions by linear functions. A differential gives a recipe for carrying out such a linear approximation. Most linear approximations can be improved by adding more terms in a Taylor series expansion. Adding quadratic terms brings in second differentials. We will meet these in the next chapter after we have mastered first differentials. Our current treatment stresses theory and counterexamples rather than the nuts and bolts of differentiation.

4.2 Univariate Derivatives

In this section we explore univariate differentiation in more detail. The standard repertoire of differentiable functions includes the derivatives nx^{n-1}

TABLE 4.1. Derivatives of some elementary functions

$f(x)$	$f'(x)$	$f(x)$	$f'(x)$	$f(x)$	$f'(x)$
x^n	nx^{n-1}	e^x	e^x	$\ln x$	$1/x$
$\sin x$	$\cos x$	$\cos x$	$-\sin x$	$\tan x$	$1 + \tan^2 x$
$\sinh x$	$\cosh x$	$\cosh x$	$\sinh x$	$\tanh x$	$1 - \tanh^2 x$
$\arcsin x$	$1/\sqrt{1-x^2}$	$\arccos x$	$-1/\sqrt{1-x^2}$	$\arctan x$	$1/(1+x^2)$
$\operatorname{arcsinh} x$	$1/\sqrt{x^2+1}$	$\operatorname{arccosh} x$	$1/\sqrt{x^2-1}$	$\operatorname{arctanh} x$	$1/(1-x^2)$

of the monomials x^n and, via the sum, product, and quotient rules, the derivatives of all polynomials and rational functions. These functions are supplemented by special functions such as $\ln x$, e^x , $\sin x$, and $\cos x$. Virtually all of the special functions can be defined by power series or as the solutions of differential equations. For instance, the system of differential equations

$$\begin{aligned}(\cos x)' &= -\sin x \\(\sin x)' &= \cos x\end{aligned}$$

with the initial conditions $\cos 0 = 1$ and $\sin 0 = 0$ determines these trigonometric functions. We will take most of these facts for granted except to add in the case of $\cos x$ and $\sin x$ that the solution of the defining system of differential equations involves a particular matrix exponential. Table 4.1 lists the derivatives of the most important elementary functions.

It is worth emphasizing that differentiation is a purely local operation and that differentiability at a point implies continuity at the same point. The converse is clearly false. The functions

$$f_n(x) = \begin{cases} x^n & x \text{ rational} \\ 0 & x \text{ irrational} \end{cases}$$

illustrate the local character of continuity and differentiability. For $n > 0$ the functions $f_n(x)$ are continuous at the point 0 but discontinuous everywhere else. In contrast, $f'_1(0)$ fails to exist while $f'_n(0) = 0$ for all $n \geq 2$. In this instance, we must resort directly to the definition (3.1) to evaluate derivatives.

We have already mentioned Fermat's result that $f'(x)$ must vanish at any interior extreme point. For example, suppose that c is a local maximum of $f(x)$ on (a, b) . If $f'(c) > 0$, then choose $\epsilon > 0$ such that $f'(c) - \epsilon > 0$. This choice then entails

$$\begin{aligned}f(x) &> f(c) + [f'(c) - \epsilon](x - c) \\ &> f(c)\end{aligned}$$

for all $x > c$ with $x - c$ sufficiently small, contradicting the assumption that c is a local maximum. If $f'(c) < 0$, we reach a similar contradiction using nearby points on the left of c .

Fermat's principle has some surprising implications. Among these is the mean value theorem.

Proposition 4.2.1 *Suppose $f(x)$ is continuous on $[a, b]$ and differentiable on (a, b) . Then there exists a point $c \in (a, b)$ such that*

$$f(b) - f(a) = f'(c)(b - a).$$

Proof: Consider the function

$$g(x) = f(b) - f(x) + \frac{f(b) - f(a)}{b - a}(x - b).$$

Clearly, $g(x)$ is also continuous on $[a, b]$ and differentiable on (a, b) . Furthermore, $g(a) = g(b) = 0$. It follows that $g(x)$ attains either a maximum or a minimum at some $c \in (a, b)$. At this point, $g'(c) = 0$, which is equivalent to the mean value property. ■

The mean value theorem has the following consequences:

- (a) If $f'(x) \geq 0$ for all $x \in (a, b)$, then $f(x)$ is increasing.
- (b) If $f'(x) = 0$ for all $x \in (a, b)$, then $f(x)$ is constant.
- (c) If $f'(x) \leq 0$ for all $x \in (a, b)$, then $f(x)$ is decreasing.

For an alternative proof, one can build on Example 3.2.4 of Chap. 3. See Problem 4 of that chapter.

Example 4.2.1 *A Trigonometric Identity*

The function $f(x) = \cos^2 x + \sin^2 x$ has derivative

$$f'(x) = -2 \cos x \sin x + 2 \sin x \cos x = 0.$$

Therefore, $f(x) = f(0) = 1$ for all x . ■

Here are two related matrix applications of univariate differentiation.

Example 4.2.2 *Differential Equations and the Matrix Exponential*

The derivative of a vector or matrix-valued function $f(x)$ with domain (a, b) is defined entry by entry. We have already met the matrix-valued differential equation $N'(t) = \mathbf{M}N(t)$ with initial condition $N(0) = \mathbf{I}$. To demonstrate that $N(t) = e^{t\mathbf{M}}$ is a solution, consider the difference quotient

$$\begin{aligned} \frac{e^{(t+s)\mathbf{M}} - e^{t\mathbf{M}}}{s} &= \frac{1}{s} \sum_{j=1}^{\infty} \frac{(t+s)^j - t^j}{j!} \mathbf{M}^j \\ &= \mathbf{M} \sum_{j=1}^{\infty} \frac{t^{j-1}}{(j-1)!} \mathbf{M}^{j-1} + \sum_{j=1}^{\infty} \frac{(t+s)^j - t^j - jt^{j-1}s}{sj!} \mathbf{M}^j \\ &= \mathbf{M}e^{t\mathbf{M}} + \sum_{j=1}^{\infty} \frac{(t+s)^j - t^j - jt^{j-1}s}{sj!} \mathbf{M}^j. \end{aligned}$$

We now apply the error estimate (3.10) of Chap. 3 for the first-order Taylor expansion of the function $f(t) = t^j$. If c bounds $|t+s|$ for s near 0, it follows that

$$|(t+s)^j - t^j - jt^{j-1}s| \leq j(j-1)c^{j-2}s^2/2$$

and that

$$\begin{aligned} \left\| \sum_{j=1}^{\infty} \frac{(t+s)^j - t^j - jt^{j-1}s}{sj!} \mathbf{M}^j \right\| &\leq \frac{|s|}{2} \sum_{j=2}^{\infty} \frac{j(j-1)c^{j-2}}{j!} \|\mathbf{M}\|^j \\ &= \frac{|s|}{2} \|\mathbf{M}\|^2 e^{c\|\mathbf{M}\|}. \end{aligned}$$

This is enough to show that

$$\lim_{s \rightarrow 0} \left\| \frac{e^{(t+s)\mathbf{M}} - e^{t\mathbf{M}}}{s} - \mathbf{M}e^{t\mathbf{M}} \right\| = 0.$$

One can demonstrate that $e^{t\mathbf{M}}$ is the unique solution of the differential equation $N'(t) = \mathbf{M}N(t)$ subject to $N(0) = \mathbf{I}$ by considering the matrix $P(t) = e^{-t\mathbf{M}}N(t)$ using any solution $N(t)$. Because the product rule of differentiation pertains to matrix multiplication as well as to ordinary multiplication,

$$P'(t) = -\mathbf{M}e^{-t\mathbf{M}}N(t) + e^{-t\mathbf{M}}\mathbf{M}N(t) = \mathbf{0}.$$

By virtue of part (b) of Proposition 4.2.1, $P(t)$ is the constant matrix $P(0) = \mathbf{I}$. If we take $N(t) = e^{t\mathbf{M}}$, then this argument demonstrates that $e^{-t\mathbf{M}}$ is the matrix inverse of $e^{t\mathbf{M}}$. If we take $N(t)$ to be an arbitrary solution of the differential equation, then multiplying both sides of $e^{-t\mathbf{M}}N(t) = \mathbf{I}$ on the left by $e^{t\mathbf{M}}$ implies that $N(t) = e^{t\mathbf{M}}$ as claimed. ■

Example 4.2.3 Matrix Logarithm

Let \mathbf{M} be a square matrix with $\|\mathbf{M}\| < 1$. It is tempting to define the logarithm of $\mathbf{I} - \mathbf{M}$ by the series expansion

$$\ln(\mathbf{I} - \mathbf{M}) = -\sum_{k=1}^{\infty} \frac{\mathbf{M}^k}{k}$$

valid for scalars. This definition does not settle the question of whether

$$e^{\ln(\mathbf{I} - \mathbf{M})} = \mathbf{I} - \mathbf{M}. \quad (4.1)$$

The traditional approach to such issues relies on Jordan canonical forms [137]. Here we would like to sketch an analytic proof. Consider the matrix-valued functions

$$\begin{aligned} f(t) &= e^{\ln(\mathbf{I} - t\mathbf{M})} \\ f_n(t) &= e^{-\sum_{k=1}^n t^k \mathbf{M}^k / k} \\ &= e^{-t\mathbf{M}} e^{-t^2 \mathbf{M}^2 / 2} \dots e^{-t^n \mathbf{M}^n / n} \end{aligned}$$

of the scalar t . It is clear that $f_n(t)$ converges uniformly to $f(t)$ on every interval $[0, 1 + \delta)$ for $\delta > 0$ small enough. Furthermore, the product rule, the chain rule, and the law of exponents show that

$$\begin{aligned} f'_n(t) &= -(\mathbf{M} + t\mathbf{M}^2 + \cdots + t^{n-1}\mathbf{M}^n)f_n(t) \\ &= -\mathbf{M}(\mathbf{I} - t^n\mathbf{M}^n)(\mathbf{I} - t\mathbf{M})^{-1}f_n(t). \end{aligned}$$

Because $\mathbf{I} - t^n\mathbf{M}^n$ tends to \mathbf{I} , it follows that

$$\lim_{n \rightarrow \infty} f'_n(t) = -\mathbf{M}(\mathbf{I} - t\mathbf{M})^{-1}f(t)$$

uniformly on $[0, 1 + \delta)$. This in turn implies

$$\begin{aligned} f(t) - f(0) &= \lim_{n \rightarrow \infty} [f_n(t) - f_n(0)] \\ &= \lim_{n \rightarrow \infty} \int_0^t f'_n(s) ds \\ &= -\mathbf{M} \int_0^t (\mathbf{I} - s\mathbf{M})^{-1} f(s) ds \end{aligned}$$

by virtue of Problem 17 of Chap. 3. Differentiating this equation with respect to t produces the differential equation

$$f'(t) = -\mathbf{M}(\mathbf{I} - t\mathbf{M})^{-1}f(t) \quad (4.2)$$

with initial condition $f(0) = \mathbf{I}$. Clearly $f(t) = \mathbf{I} - t\mathbf{M}$ is one solution of the differential equation (4.2). In view of Problem 16, this solution is unique. Comparing the two formulas for $f(t)$ at the point $t = 1$ now gives the desired conclusion (4.1). ■

4.3 Partial Derivatives

There are several possible ways to extend differentiation to real-valued functions on \mathbb{R}^n . The most familiar is the partial derivative

$$\partial_i f(\mathbf{x}) = \frac{\partial}{\partial x_i} f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{e}_i) - f(\mathbf{x})}{t},$$

where \mathbf{e}_i is one of the standard unit vectors spanning \mathbb{R}^n . There is nothing sacred about the coordinate directions. The directional derivative along the direction \mathbf{v} is

$$d_{\mathbf{v}} f(\mathbf{x}) = \lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}.$$

If we confine $t \geq 0$ in this limit, then we have a forward derivative along \mathbf{v} . Readers should be on the alert that we will use the same symbol $d_{\mathbf{v}}f(\mathbf{x})$ for both directional derivatives. The context will make it clear which concept is pertinent.

To illustrate these definitions, consider the function

$$f(\mathbf{x}) = \sqrt{|x_1 x_2|}$$

on \mathbb{R}^2 . It is clear that both partial derivatives are 0 at the origin. Along a direction $\mathbf{v} = (v_1, v_2)^*$ with neither $v_1 = 0$ nor $v_2 = 0$, consider the difference quotient

$$\frac{f(\mathbf{0} + t\mathbf{v}) - f(\mathbf{0})}{t} = \frac{|t|\sqrt{|v_1 v_2|}}{t}.$$

This has limit $\sqrt{|v_1 v_2|}$ as long as we restrict $t > 0$. Thus, the forward directional derivative exists, but the full directional derivative does not.

For another example, let

$$f(\mathbf{x}) = \begin{cases} x_1 + x_2 & \text{if } x_1 = 0 \text{ or } x_2 = 0 \\ 1 & \text{otherwise.} \end{cases}$$

This function is clearly discontinuous at the origin of \mathbb{R}^2 , but the partial derivatives $\partial_1 f(\mathbf{0}) = 1$ and $\partial_2 f(\mathbf{0}) = 1$ are well defined there. These and similar anomalies suggest the need for a carefully structured theory of differentiability. Such a theory is presented in the next section.

Second and higher-order partial derivatives are defined in the obvious way. For typographical convenience, we will occasionally employ such abbreviations as

$$\partial_{ij}^2 f(\mathbf{x}) = \frac{\partial^2}{\partial x_i \partial x_j} f(\mathbf{x}).$$

Readers will doubtless recall from calculus the equality of mixed second partial derivatives. This property can fail. For example, suppose we define $f(\mathbf{x}) = g(x_1)$, where $g(x_1)$ is nowhere differentiable. Then $\partial_2 f(\mathbf{x})$ and $\partial_{12}^2 f(\mathbf{x})$ are identically 0 while $\partial_1 f(\mathbf{x})$ does not even exist. The key to restoring harmony is to impose continuity in a neighborhood of the current point.

Proposition 4.3.1 *Suppose the real-valued function $f(\mathbf{y})$ on \mathbb{R}^2 has partial derivatives $\partial_1 f(\mathbf{y})$, $\partial_2 f(\mathbf{y})$, and $\partial_{12}^2 f(\mathbf{y})$ on some open set. If $\partial_{12}^2 f(\mathbf{y})$ is continuous at a point \mathbf{x} in the set, then $\partial_{21}^2 f(\mathbf{x})$ exists and*

$$\frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x}) = \partial_{21}^2 f(\mathbf{x}) = \partial_{12}^2 f(\mathbf{x}) = \frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x}). \quad (4.3)$$

This result extends in the obvious way to the equality of second mixed partials for functions defined on open subsets of \mathbb{R}^n for $n > 2$.

Proof: Consider the first difference by u_1

$$g(x_2) = \Delta_1 f(x_1, x_2) = f(x_1 + u_1, x_2) - f(x_1, x_2)$$

and the second difference by u_2

$$\begin{aligned} \Delta_{21} f(x_1, x_2) &= g(x_2 + u_2) - g(x_2) \\ &= f(x_1 + u_1, x_2 + u_2) - f(x_1, x_2 + u_2) \\ &\quad - f(x_1 + u_1, x_2) + f(x_1, x_2). \end{aligned}$$

Applying the mean value theorem twice gives

$$\begin{aligned} \Delta_{21} f(x_1, x_2) &= u_2 g'(x_2 + \theta_2 u_2) \\ &= u_2 [\partial_2 f(x_1 + u_1, x_2 + \theta_2 u_2) - \partial_2 f(x_1, x_2 + \theta_2 u_2)] \\ &= u_1 u_2 \partial_{12}^2 f(x_1 + \theta_1 u_1, x_2 + \theta_2 u_2) \end{aligned}$$

for θ_1 and θ_2 in $(0, 1)$. In view of the continuity of $\partial_{12}^2 f(\mathbf{y})$ at \mathbf{x} , it follows that

$$\lim_{\|\mathbf{u}\| \rightarrow 0} \frac{\Delta_{21} f(x_1, x_2)}{u_1 u_2} = \partial_{12}^2 f(x_1, x_2)$$

regardless of how the limit is approached. This proves the existence of the iterated limit

$$\begin{aligned} \lim_{u_2 \rightarrow 0} \lim_{u_1 \rightarrow 0} \frac{\Delta_{21} f(x_1, x_2)}{u_1 u_2} &= \lim_{u_2 \rightarrow 0} \frac{\partial_1 f(x_1, x_2 + u_2) - \partial_1 f(x_1, x_2)}{u_2} \\ &= \partial_{21}^2 f(x_1, x_2) \end{aligned}$$

and simultaneously the equality of mixed partials. Problem 22 of Chap. 3 offers an alternative proof under stronger hypotheses. ■

4.4 Differentials

The question of when a real-valued function is differentiable is perplexing because of the variety of possible definitions. In choosing an appropriate definition, we are governed by several considerations. First, it should be consistent with the classical definition of differentiability on the real line. Second, continuity at a point should be a consequence of differentiability at the point. Third, all directional derivatives should exist. Fourth, the differential should vanish wherever the function attains a local maximum or minimum on the interior of its domain. Fifth, the standard rules for combining differentiable functions should apply. Sixth, the logical proofs of

the rules should be as transparent as possible. Seventh, the extension to vector-valued functions should be painless. Eighth and finally, our geometric intuition should be enhanced.

We now present a definition conceived by Constantin Carathéodory [40] and expanded by recent authors [1, 18, 29, 160] that fulfills these conditions. A real-valued function $f(\mathbf{y})$ on an open set $S \subset \mathbb{R}^m$ is said to be differentiable at $\mathbf{x} \in S$ if a function $s(\mathbf{y}, \mathbf{x})$ exists for \mathbf{y} near \mathbf{x} satisfying

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}) \\ \lim_{\mathbf{y} \rightarrow \mathbf{x}} s(\mathbf{y}, \mathbf{x}) &= s(\mathbf{x}, \mathbf{x}). \end{aligned} \quad (4.4)$$

The row vector $s(\mathbf{y}, \mathbf{x})$ is called a slope function. We will see in a moment that its limit $s(\mathbf{x}, \mathbf{x})$ defines the differential $df(\mathbf{x})$ of $f(\mathbf{y})$ at \mathbf{x} .

The standard definition of differentiability due to Fréchet reads

$$f(\mathbf{y}) - f(\mathbf{x}) = df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + o(\|\mathbf{y} - \mathbf{x}\|)$$

for \mathbf{y} near the point \mathbf{x} . The row vector $df(\mathbf{x})$ appearing here is again termed the differential of $f(\mathbf{y})$ at \mathbf{x} . Fréchet's definition is less convenient than Carathéodory's because the former invokes approximate equality rather than true equality. Observe that the error $[s(\mathbf{y}, \mathbf{x}) - df(\mathbf{x})](\mathbf{y} - \mathbf{x})$ under Carathéodory's definition satisfies

$$|[s(\mathbf{y}, \mathbf{x}) - df(\mathbf{x})](\mathbf{y} - \mathbf{x})| \leq \|s(\mathbf{y}, \mathbf{x}) - df(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\|,$$

which is $o(\|\mathbf{y} - \mathbf{x}\|)$ as \mathbf{y} tends to \mathbf{x} in view of the continuity of $s(\mathbf{y}, \mathbf{x})$. Thus, Carathéodory's definition implies Fréchet's definition. The converse is trivial when the argument of $f(x)$ is a scalar, for then the difference quotient

$$s(y, x) = \frac{f(y) - f(x)}{y - x} \quad (4.5)$$

serves as a slope function. Proposition 4.4.1 addresses the general case.

Carathéodory's definition (4.4) has some immediate consequences. For example, it obviously compels $f(\mathbf{y})$ to be continuous at \mathbf{x} . Furthermore, if we send t to 0 in the equation

$$\frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = s(\mathbf{x} + t\mathbf{v}, \mathbf{x})\mathbf{v}, \quad (4.6)$$

then it is clear that the directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ exists and equals $s(\mathbf{x}, \mathbf{x})\mathbf{v}$. The special case $\mathbf{v} = \mathbf{e}_i$ shows that the i th component of $s(\mathbf{x}, \mathbf{x})$ reduces to the partial derivative $\partial_i f(\mathbf{x})$. Since $s(\mathbf{x}, \mathbf{x})$ and $df(\mathbf{x})$ agree component by component, they are equal, and, in general, we have the formula $d_{\mathbf{v}}f(\mathbf{x}) = df(\mathbf{x})\mathbf{v}$ for the directional derivative.

Fermat's principle that the differential of a function vanishes at an interior maximum or minimum point is also trivial to check in this context. Suppose \mathbf{x} affords a local minimum of $f(\mathbf{y})$. Then

$$f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x}) = ts(\mathbf{x} + t\mathbf{v}, \mathbf{x})\mathbf{v} \geq 0$$

for all \mathbf{v} and small $t > 0$. Taking limits in the identity (4.6) now yields the conclusion $df(\mathbf{x})\mathbf{v} \geq 0$. The only way this can hold for all \mathbf{v} is for $df(\mathbf{x}) = \mathbf{0}$. If \mathbf{x} occurs on the boundary of the domain of $f(\mathbf{y})$, then we can still glean useful information. For example, if $f(y)$ is differentiable on the closed interval $[c, d]$ and c provides a local minimum, then the condition $f'(c) \geq 0$ must hold.

The extension of the definition of differentiability to vector-valued functions is equally simple. Suppose $f(\mathbf{y})$ maps an open subset $S \subset \mathbb{R}^m$ into \mathbb{R}^n . Then $f(\mathbf{y})$ is said to be differentiable at $\mathbf{x} \in S$ if there exists an $n \times m$ matrix-valued function $s(\mathbf{y}, \mathbf{x})$ continuous at \mathbf{x} and satisfying equation (4.4) for \mathbf{y} near \mathbf{x} . The limit $\lim_{\mathbf{y} \rightarrow \mathbf{x}} s(\mathbf{y}, \mathbf{x}) = df(\mathbf{x})$ is again called the differential of $f(\mathbf{y})$ at \mathbf{x} . The rows of the differential are the differentials of the component functions of $f(\mathbf{x})$. Thus, $f(\mathbf{y})$ is differentiable at \mathbf{x} if and only if each of its components is differentiable at \mathbf{x} . This characterization is also valid under Fréchet's definition of the differential and leads to a simple proof of the second half of the next proposition.

Proposition 4.4.1 *Carathéodory's definition and Fréchet's definition of the differential are logically equivalent.*

Proof: We have already proved that Carathéodory's definition implies Fréchet's definition. The converse is valid because it is valid for scalar-valued functions. For a matrix-oriented proof of the converse, suppose that $f(\mathbf{y})$ is Fréchet differentiable at \mathbf{x} . If we define the slope function

$$s(\mathbf{y}, \mathbf{x}) = \frac{1}{\|\mathbf{y} - \mathbf{x}\|^2} [f(\mathbf{y}) - f(\mathbf{x}) - df(\mathbf{x})(\mathbf{y} - \mathbf{x})](\mathbf{y} - \mathbf{x})^* + df(\mathbf{x})$$

for $\mathbf{y} \neq \mathbf{x}$, then the identity $f(\mathbf{y}) - f(\mathbf{x}) = s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$ certainly holds. To show that $s(\mathbf{y}, \mathbf{x})$ tends to $df(\mathbf{x})$ as \mathbf{y} tends to \mathbf{x} , we now observe that $s(\mathbf{y}, \mathbf{x}) = \mathbf{u}\mathbf{v}^* + df(\mathbf{x})$ for vectors \mathbf{u} and \mathbf{v} . In view of the Cauchy-Schwarz inequality, the spectral norm of the matrix outer product $\mathbf{u}\mathbf{v}^*$ satisfies

$$\|\mathbf{u}\mathbf{v}^*\| = \sup_{\mathbf{w} \neq \mathbf{0}} \frac{\|\mathbf{u}\mathbf{v}^*\mathbf{w}\|}{\|\mathbf{w}\|} = \sup_{\mathbf{w} \neq \mathbf{0}} \frac{|\mathbf{v}^*\mathbf{w}|\|\mathbf{u}\|}{\|\mathbf{w}\|} \leq \|\mathbf{u}\|\|\mathbf{v}\|.$$

In the current setting this translates into the inequality

$$\|\mathbf{u}\mathbf{v}^*\| \leq \frac{\|f(\mathbf{y}) - f(\mathbf{x}) - df(\mathbf{x})(\mathbf{y} - \mathbf{x})\|}{\|\mathbf{y} - \mathbf{x}\|} \cdot \frac{\|\mathbf{y} - \mathbf{x}\|}{\|\mathbf{y} - \mathbf{x}\|}.$$

Hence, Fréchet's condition implies $\lim_{\mathbf{y} \rightarrow \mathbf{x}} \|s(\mathbf{y}, \mathbf{x}) - df(\mathbf{x})\| = 0$. ■

In many cases it is easy to identify a slope function. We have already mentioned slope functions defined by difference quotients. It is worth stressing that while differentials are uniquely determined, slope functions are not. The real-valued function $f(\mathbf{y}) = y_1y_2$ is typical. Indeed, the identities

$$\begin{aligned}y_1y_2 - x_1x_2 &= x_2(y_1 - x_1) + y_1(y_2 - x_2) \\y_1y_2 - x_1x_2 &= y_2(y_1 - x_1) + x_1(y_2 - x_2)\end{aligned}$$

define two equally valid slope functions at \mathbf{x} .

Example 4.4.1 *Differentials of Linear and Quadratic Functions*

A linear transformation $f(\mathbf{y}) = \mathbf{M}\mathbf{y}$ is differentiable with slope function $s(\mathbf{y}, \mathbf{x}) = \mathbf{M}$. The real-valued coordinate functions y_i of $\mathbf{y} \in \mathbb{R}^n$ fall into this category. For a symmetric matrix \mathbf{M} , the quadratic form $g(\mathbf{x}) = \mathbf{x}^*\mathbf{M}\mathbf{x}$ has the difference

$$\mathbf{y}^*\mathbf{M}\mathbf{y} - \mathbf{x}^*\mathbf{M}\mathbf{x} = (\mathbf{y} + \mathbf{x})^*\mathbf{M}(\mathbf{y} - \mathbf{x}).$$

This gives the differential $dg(\mathbf{x}) = 2\mathbf{x}^*\mathbf{M}$ and gradient $\nabla g(\mathbf{x}) = 2\mathbf{M}\mathbf{x}$. ■

Example 4.4.2 *Differential of a Multilinear Map*

A multilinear map $M[\mathbf{u}_1, \dots, \mathbf{u}_k]$ as defined in Example 2.5.10 is differentiable. The expansion

$$\begin{aligned}M[\mathbf{v}_1, \dots, \mathbf{v}_k] - M[\mathbf{u}_1, \dots, \mathbf{u}_k] &= M[\mathbf{v}_1 - \mathbf{u}_1, \mathbf{v}_2, \dots, \mathbf{v}_k] \\ &+ M[\mathbf{u}_1, \mathbf{v}_2 - \mathbf{u}_2, \dots, \mathbf{v}_k] \\ &\vdots \\ &+ M[\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{v}_k - \mathbf{u}_k]\end{aligned}$$

displays the slope function as a sum. The corresponding differential

$$\begin{aligned}dM[\mathbf{u}_1, \dots, \mathbf{u}_k][\mathbf{w}_1, \dots, \mathbf{w}_k] &= M[\mathbf{w}_1, \mathbf{u}_2, \dots, \mathbf{u}_k] \\ &+ M[\mathbf{u}_1, \mathbf{w}_2, \dots, \mathbf{u}_k] \\ &\vdots \\ &+ M[\mathbf{u}_1, \mathbf{u}_2, \dots, \mathbf{w}_k].\end{aligned}$$

emerges in the limit. ■

The rules for calculating differentials of algebraic combinations of differentiable functions flow easily from Carathéodory's definition.

Proposition 4.4.2 *If the two functions $f(\mathbf{x})$ and $g(\mathbf{x})$ map the open set $S \subset \mathbb{R}^m$ differentiably into \mathbb{R}^n , then the following combinations are differentiable and have the displayed differentials:*

(a) $d[\alpha f(\mathbf{x}) + \beta g(\mathbf{x})] = \alpha df(\mathbf{x}) + \beta dg(\mathbf{x})$ for all constants α and β .

(b) $d[f(\mathbf{x})^*g(\mathbf{x})] = f(\mathbf{x})^*dg(\mathbf{x}) + g(\mathbf{x})^*df(\mathbf{x})$.

(c) $d[f(\mathbf{x})^{-1}] = -f(\mathbf{x})^{-2}df(\mathbf{x})$ when $n = 1$ and $f(\mathbf{x}) \neq 0$.

Proof: Let the slope functions for $f(\mathbf{x})$ and $g(\mathbf{x})$ at \mathbf{x} be $s_f(\mathbf{y}, \mathbf{x})$ and $s_g(\mathbf{y}, \mathbf{x})$. Rule (a) follows by taking the limit of the slope function identified in the equality

$$\alpha f(\mathbf{y}) + \beta g(\mathbf{y}) - \alpha f(\mathbf{x}) - \beta g(\mathbf{x}) = [\alpha s_f(\mathbf{y}, \mathbf{x}) + \beta s_g(\mathbf{y}, \mathbf{x})](\mathbf{y} - \mathbf{x}).$$

Rule (b) stems from the equality

$$\begin{aligned} & f(\mathbf{y})^*g(\mathbf{y}) - f(\mathbf{x})^*g(\mathbf{x}) \\ &= [f(\mathbf{y}) - f(\mathbf{x})]^*g(\mathbf{y}) + f(\mathbf{x})^*[g(\mathbf{y}) - g(\mathbf{x})] \\ &= g(\mathbf{y})^*s_f(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}) + f(\mathbf{x})^*s_g(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}), \end{aligned}$$

and rule (c) stems from the equality

$$\begin{aligned} f(\mathbf{y})^{-1} - f(\mathbf{x})^{-1} &= -f(\mathbf{y})^{-1}f(\mathbf{x})^{-1}[f(\mathbf{y}) - f(\mathbf{x})] \\ &= -f(\mathbf{y})^{-1}f(\mathbf{x})^{-1}s_f(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}). \end{aligned}$$

■

The chain rule also has an beautifully straightforward proof.

Proposition 4.4.3 *Suppose $f(\mathbf{x})$ maps the open set $S \subset \mathbb{R}^k$ differentiably into \mathbb{R}^m and $g(\mathbf{z})$ maps the open set $T \subset \mathbb{R}^m$ differentiably into \mathbb{R}^n . If the image $f(S)$ is contained in T , then the composition $g \circ f(\mathbf{x})$ is differentiable with differential $dg \circ f(\mathbf{x})df(\mathbf{x})$.*

Proof: Let $s_f(\mathbf{y}, \mathbf{x})$ be the slope function of $f(\mathbf{y})$ for \mathbf{y} near \mathbf{x} and $s_g(\mathbf{z}, \mathbf{w})$ be the slope function of $g(\mathbf{z})$ for \mathbf{z} near $\mathbf{w} = f(\mathbf{x})$. The chain rule follows after taking the limit of the slope function identified in the equality

$$\begin{aligned} g \circ f(\mathbf{y}) - g \circ f(\mathbf{x}) &= s_g[f(\mathbf{y}), f(\mathbf{x})][f(\mathbf{y}) - f(\mathbf{x})] \\ &= s_g[f(\mathbf{y}), f(\mathbf{x})]s_f(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}). \end{aligned}$$

The chain rule is much harder to prove under Fréchet's definition. ■

Of course, these rules do not exhaust the techniques for finding differentials. Here is an example where we must fall back on the definition. See Problem 17 for a generalization.

Example 4.4.3 *Differential of $f_+(\mathbf{x})^2$*

Suppose $f(\mathbf{x})$ is a real-valued differentiable function. Define the function $f_+(\mathbf{x}) = \max\{f(\mathbf{x}), 0\}$. In general, $f_+(\mathbf{x})$ is not a differentiable function,

but $g(\mathbf{x}) = f_+(\mathbf{x})^2$ is. This is obvious on the open set $\{\mathbf{x} : f(\mathbf{x}) < 0\}$, where $dg(\mathbf{x}) = \mathbf{0}$, and on the open set $\{\mathbf{x} : f(\mathbf{x}) > 0\}$, where

$$dg(\mathbf{x}) = 2f(\mathbf{x})df(\mathbf{x}).$$

The troublesome points are those with $f(\mathbf{x}) = 0$. Near such a point we have $f(\mathbf{y}) - 0 = s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$, and

$$g(\mathbf{y}) - 0 = f_+(\mathbf{y})s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}).$$

It follows that $dg(\mathbf{x}) = f_+(\mathbf{x})df(\mathbf{x}) = \mathbf{0}^*$ when $f(\mathbf{x}) = 0$. In general, all three cases can be summarized by the same rule $dg(\mathbf{x}) = 2f_+(\mathbf{x})df(\mathbf{x})$. ■

Example 4.4.4 *Forward Directional Derivative of $\max_{1 \leq i \leq p} g_i(\mathbf{x})$*

As the example $|x| = \max\{-x, x\}$ illustrates, the maximum of two differentiable functions may not be differentiable. For many purposes in optimization, forward directional derivatives are adequate. Consider the maximum $f(\mathbf{x}) = \max_{1 \leq i \leq p} g_i(\mathbf{x})$ of a finite number of real-valued functions differentiable at the point \mathbf{y} . To show that $f(\mathbf{x})$ possesses all possible forward directional derivatives at \mathbf{y} , let $\mathbf{v} \neq \mathbf{0}$ be an arbitrary direction and t_n any sequence of positive numbers converging to 0. It suffices to prove that the difference quotients $t_n^{-1}[f(\mathbf{y} + t_n\mathbf{v}) - f(\mathbf{y})]$ tend to a limit $d_{\mathbf{v}}f(\mathbf{y})$ independent of the specific sequence t_n . Because the $g_i(\mathbf{x})$ are differentiable at \mathbf{y} , they are also continuous at \mathbf{y} . Those $g_i(\mathbf{x})$ with $g_i(\mathbf{y}) < f(\mathbf{y})$ play no role in determining $f(\mathbf{x})$ near \mathbf{y} and can be discarded in calculating a directional derivative. Hence, we assume without loss of generality that all $g_i(\mathbf{y}) = f(\mathbf{y})$. With this proviso, we claim that $d_{\mathbf{v}}f(\mathbf{y}) = \max_{1 \leq i \leq p} dg_i(\mathbf{y})\mathbf{v}$. The inequality

$$\liminf_{n \rightarrow \infty} \frac{f(\mathbf{y} + t_n\mathbf{v}) - f(\mathbf{y})}{t_n} \geq \liminf_{n \rightarrow \infty} \frac{g_i(\mathbf{y} + t_n\mathbf{v}) - g_i(\mathbf{y})}{t_n} = dg_i(\mathbf{y})\mathbf{v}$$

for any i is obvious. Suppose that

$$\limsup_{n \rightarrow \infty} \frac{f(\mathbf{y} + t_n\mathbf{v}) - f(\mathbf{y})}{t_n} > \max_{1 \leq i \leq p} dg_i(\mathbf{y})\mathbf{v}. \quad (4.7)$$

In view of the definition of \limsup , there exists an $\epsilon > 0$ and a subsequence t_{n_m} along which

$$\frac{f(\mathbf{y} + t_{n_m}\mathbf{v}) - f(\mathbf{y})}{t_{n_m}} \geq \max_{1 \leq i \leq p} dg_i(\mathbf{y})\mathbf{v} + \epsilon.$$

Passing to a subsubsequence if necessary, we can choose a j such that

$$\frac{f(\mathbf{y} + t_{n_m}\mathbf{v}) - f(\mathbf{y})}{t_{n_m}} = \frac{g_j(\mathbf{y} + t_{n_m}\mathbf{v}) - g_j(\mathbf{y})}{t_{n_m}}$$

for all m . Taking limits now produces the contradiction

$$dg_j(\mathbf{y})\mathbf{v} \geq \max_{1 \leq i \leq p} dg_i(\mathbf{y})\mathbf{v} + \epsilon.$$

Hence, inequality (4.7) is false, and the difference quotients tend to the claimed limit. Appendix A.6 treats this example in more depth. ■

Example 4.4.5 Analytic Functions

When we come to complex-valued functions $f(\mathbf{z})$ of a complex variable \mathbf{z} , we are confronted with a dilemma in defining differentiability. On the one hand, we can substitute complex arithmetic operations for real arithmetic operations in Carathéodory's definition. If we adopt this perspective, then the equations

$$\begin{aligned} f(\mathbf{z}) - f(\mathbf{w}) &= s(\mathbf{z}, \mathbf{w})(\mathbf{z} - \mathbf{w}) \\ \lim_{\mathbf{z} \rightarrow \mathbf{w}} s(\mathbf{z}, \mathbf{w}) &= s(\mathbf{w}, \mathbf{w}) \end{aligned} \quad (4.8)$$

are summarized by saying that $f(\mathbf{z})$ is analytic (or holomorphic) at \mathbf{w} . Most of the results we have proved for differentiable functions carry over without change to analytic functions. The chain rule is a case in point. On the other hand, we can view the complex plane as \mathbb{R}^2 and decompose $\mathbf{z} = x + iy$ and $f(\mathbf{z}) = g(\mathbf{z}) + ih(\mathbf{z})$ into their real and imaginary parts with $i = \sqrt{-1}$. In this context, $f(\mathbf{z})$ is differentiable at \mathbf{w} if there exists a 2×2 slope matrix $m(\mathbf{z}, \mathbf{w})$ satisfying

$$\begin{aligned} \begin{pmatrix} g(\mathbf{z}) - g(\mathbf{w}) \\ h(\mathbf{z}) - h(\mathbf{w}) \end{pmatrix} &= m(\mathbf{z}, \mathbf{w}) \begin{pmatrix} x - u \\ y - v \end{pmatrix} \\ \lim_{\mathbf{z} \rightarrow \mathbf{w}} m(\mathbf{z}, \mathbf{w}) &= m(\mathbf{w}, \mathbf{w}) \end{aligned} \quad (4.9)$$

for $\mathbf{w} = u + iv$. A function $f(\mathbf{z})$ analytic at \mathbf{w} is differentiable at \mathbf{w} . Indeed, if we decompose $s(\mathbf{z}, \mathbf{w}) = r(\mathbf{z}, \mathbf{w}) + it(\mathbf{z}, \mathbf{w})$ in equation (4.8), then identifying real and imaginary parts demonstrates that the matrix

$$m(\mathbf{z}, \mathbf{w}) = \begin{pmatrix} r(\mathbf{z}, \mathbf{w}) & -t(\mathbf{z}, \mathbf{w}) \\ t(\mathbf{z}, \mathbf{w}) & r(\mathbf{z}, \mathbf{w}) \end{pmatrix}$$

satisfies equation (4.9). Hence, analyticity is stronger than differentiability. Taking limits on \mathbf{z} in this definition of $m(\mathbf{z}, \mathbf{w})$ furthermore shows that $f(\mathbf{z})$ satisfies the Cauchy-Riemann equations

$$\begin{aligned} \frac{\partial}{\partial x} g(\mathbf{z}) &= \frac{\partial}{\partial y} h(\mathbf{z}) \\ -\frac{\partial}{\partial y} g(\mathbf{z}) &= \frac{\partial}{\partial x} h(\mathbf{z}) \end{aligned}$$

at $\mathbf{z} = \mathbf{w}$. A beautiful theory with surprising consequences can be constructed for analytic functions [129, 244]. For example, every analytic function on an open domain is infinitely differentiable on that domain. The price we pay for such powerful results is that the class of analytic functions is much smaller than the class of differentiable functions of \mathbb{R}^2 into itself. For example the failure of the Cauchy-Riemann equations implies that the complex conjugate function $x + iy \mapsto x - iy$ is not analytic even though it is differentiable. ■

4.5 Multivariate Mean Value Theorem

The mean value theorem is one of the most useful tools of the differential calculus. Here is a simple generalization to multiple dimensions.

Proposition 4.5.1 *Let the function $f(\mathbf{y})$ map an open subset S of \mathbb{R}^m to \mathbb{R}^n . If $f(\mathbf{y})$ is differentiable on a neighborhood of $\mathbf{x} \in S$, then*

$$f(\mathbf{y}) = f(\mathbf{x}) + \int_0^1 df[\mathbf{x} + t(\mathbf{y} - \mathbf{x})] dt (\mathbf{y} - \mathbf{x}) \quad (4.10)$$

for \mathbf{y} near \mathbf{x} . If S is convex, then identity (4.10) holds for all $\mathbf{y} \in S$.

Proof: Integrating component by component, we need only consider the case $n = 1$. According to the chain rule stated in Proposition 4.4.3, the real-valued function $g(t) = f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})]$ of the scalar t has differential $dg(t) = df[\mathbf{x} + t(\mathbf{y} - \mathbf{x})](\mathbf{y} - \mathbf{x})$. Because differentials and derivatives coincide on the real line, equality (4.10) follows from the fundamental theorem of calculus applied to $g(t)$. ■

The notion of continuous differentiability is ambiguous. On the one hand, we could say that $f(\mathbf{y})$ is continuously differentiable around \mathbf{x} if it possesses a slope function $s(\mathbf{y}, \mathbf{z})$ that is jointly continuous in its two arguments. This implies the continuity of $df(\mathbf{y}) = s(\mathbf{y}, \mathbf{y})$ around \mathbf{x} . On the other hand, continuous differentiability suggests that we postulate the continuity of $df(\mathbf{y})$ to start with. In this case, equation (4.10) yields the slope function

$$s(\mathbf{y}, \mathbf{z}) = \int_0^1 df[\mathbf{z} + t(\mathbf{y} - \mathbf{z})] dt. \quad (4.11)$$

If $df(\mathbf{y})$ is continuous near \mathbf{x} , then this choice of $s(\mathbf{y}, \mathbf{z})$ is jointly continuous in its arguments. Hence, the two definitions of continuous differentiability coincide. It is noteworthy that the particular slope function (4.11) is symmetric in the sense that $s(\mathbf{y}, \mathbf{z}) = s(\mathbf{z}, \mathbf{y})$.

Example 4.5.1 *Characterization of Constant Functions*

If the differential $df(\mathbf{x})$ of a real-valued function $f(\mathbf{x})$ vanishes on an open connected set $S \subset \mathbb{R}^n$, then $f(\mathbf{x})$ is constant there. To establish this fact, let \mathbf{z} be an arbitrary point of S and define $T = \{\mathbf{x} \in S : f(\mathbf{x}) = f(\mathbf{z})\}$. Given the continuity of $f(\mathbf{x})$, it is obvious that T is closed relative to S . To prove that $T = S$, it suffices to show that T is also open relative to S . Indeed, if T and $S \setminus T$ are both nonempty open sets, then they disconnect S . Now any point $\mathbf{x} \in T$ is contained in a ball $B(\mathbf{x}, r) \subset S$. For $\mathbf{y} \in B(\mathbf{x}, r)$, formula (4.10) and the vanishing of the differential show that $f(\mathbf{y}) = f(\mathbf{x})$. Thus, T is open. ■

Example 4.5.2 *Failure of Proposition 4.2.1*

Consider the function $f(x) = (\cos x, \sin x)^*$ from \mathbb{R} to \mathbb{R}^2 . The obvious generalization of the mean value theorem stated in Proposition 4.2.1 fails because there is no $x \in (0, 2\pi)$ satisfying

$$\mathbf{0} = f(2\pi) - f(0) = \begin{pmatrix} -\sin x \\ \cos x \end{pmatrix} (2\pi - 0).$$

In this regard recall Example 4.2.1. ■

In spite of this counterexample, the bound

$$\begin{aligned} \|f(\mathbf{y}) - f(\mathbf{x})\| &\leq \left\| \int_0^1 df[\mathbf{x} + t(\mathbf{y} - \mathbf{x})] dt \right\| \cdot \|\mathbf{y} - \mathbf{x}\| \\ &\leq \sup_{t \in [0,1]} \|df[\mathbf{x} + t(\mathbf{y} - \mathbf{x})]\| \cdot \|\mathbf{y} - \mathbf{x}\| \end{aligned} \quad (4.12)$$

is often an adequate substitute for theoretical purposes for functions with a convex domain. See inequality (3.5) of Chap. 3 for a proof.

4.6 Inverse and Implicit Function Theorems

Two of the harder theorems involving differentials are the inverse and implicit function theorems. The definition of differentials through slope functions tends to make the proofs easier to understand. The current proof of the inverse function theorem also features an interesting optimization argument [1].

Proposition 4.6.1 *Let $f(\mathbf{x})$ map an open set $U \subset \mathbb{R}^n$ into \mathbb{R}^n . If $f(\mathbf{x})$ is continuously differentiable on U and the square matrix $df(\mathbf{x})$ is invertible at the point \mathbf{z} , then there exist neighborhoods V of \mathbf{z} and W of $f(\mathbf{z})$ such that the inverse function $g(\mathbf{y})$ satisfying $g \circ f(\mathbf{x}) = \mathbf{x}$ exists and maps W onto V . Furthermore, $g(\mathbf{x})$ is continuously differentiable with differential $dg(\mathbf{x}) = df[g(\mathbf{x})]^{-1}$.*

Proof: Let $f(\mathbf{x})$ have continuous slope function $s(\mathbf{y}, \mathbf{x})$. If $s(\mathbf{y}, \mathbf{x})$ is invertible as a square matrix, then the relations

$$f(\mathbf{y}) - f(\mathbf{x}) = s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$$

and

$$\mathbf{y} - \mathbf{x} = s(\mathbf{y}, \mathbf{x})^{-1}[f(\mathbf{y}) - f(\mathbf{x})]$$

are equivalent. Now suppose we know $f(\mathbf{x})$ has functional inverse $g(\mathbf{y})$. Exchanging $g(\mathbf{y})$ for \mathbf{y} and $g(\mathbf{x})$ for \mathbf{x} in the second relation above produces

$$g(\mathbf{y}) - g(\mathbf{x}) = s[g(\mathbf{y}), g(\mathbf{x})]^{-1}(\mathbf{y} - \mathbf{x}), \quad (4.13)$$

and taking limits gives the claimed differential, provided $g(\mathbf{y})$ is continuous. To prove the continuity of $g(\mathbf{y})$ and therefore the joint continuity of $s[g(\mathbf{y}), g(\mathbf{x})]^{-1}$, it suffices to show that $\|s[g(\mathbf{y}), g(\mathbf{x})]^{-1}\|$ is locally bounded. Continuity in this circumstance is then a consequence of the bound

$$\|g(\mathbf{y}) - g(\mathbf{x})\| \leq \|s[g(\mathbf{y}), g(\mathbf{x})]^{-1}\| \cdot \|\mathbf{y} - \mathbf{x}\|$$

flowing from equation (4.13). In view of these remarks, the difficult part of the proof consists in proving that $g(\mathbf{y})$ exists.

Given the continuous differentiability of $f(\mathbf{x})$, there is some neighborhood V of \mathbf{z} such that $s(\mathbf{y}, \mathbf{x})$ is invertible for all \mathbf{x} and \mathbf{y} in V . Furthermore, we can take V small enough so that the norm $\|s(\mathbf{y}, \mathbf{x})^{-1}\|$ is bounded there. On V , the equality $f(\mathbf{y}) - f(\mathbf{x}) = s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$ shows that $f(\mathbf{x})$ is one-to-one. Hence, all that remains is to show that we can shrink V so that $f(\mathbf{x})$ maps V onto an open subset W containing $f(\mathbf{z})$.

For some $r > 0$, the ball $B(\mathbf{z}, r)$ of radius r centered at \mathbf{z} is contained in V . The sphere $S(\mathbf{z}, r) = \{\mathbf{x} : \|\mathbf{x} - \mathbf{z}\| = r\}$ and the ball $B(\mathbf{z}, r)$ are disjoint and must have disjoint images under $f(\mathbf{x})$ because $f(\mathbf{x})$ is one-to-one on V . In particular, $f(\mathbf{z})$ is not contained in $f[S(\mathbf{z}, r)]$. The latter set is compact because $S(\mathbf{z}, r)$ is compact and $f(\mathbf{x})$ is continuous. Let $d > 0$ be the distance from $f(\mathbf{z})$ to $f[S(\mathbf{z}, r)]$.

We now define the set W mentioned in the statement of the proposition to be the ball $B[f(\mathbf{z}), d/2]$ and show that W is contained in the image of $B(\mathbf{z}, r)$ under $f(\mathbf{x})$. Take any $\mathbf{y} \in W = B[f(\mathbf{z}), d/2]$. The particular function $h(\mathbf{x}) = \|\mathbf{y} - f(\mathbf{x})\|^2$ is differentiable and attains its minimum on the closed ball $C(\mathbf{z}, r) = \{\mathbf{x} : \|\mathbf{x} - \mathbf{z}\| \leq r\}$. This minimum is strictly less than $(d/2)^2$ because \mathbf{z} certainly performs this well. Furthermore, the minimum cannot be reached at a point $\mathbf{u} \in S(\mathbf{z}, r)$, for then

$$\|f(\mathbf{u}) - f(\mathbf{z})\| \leq \|f(\mathbf{u}) - \mathbf{y}\| + \|\mathbf{y} - f(\mathbf{z})\| < 2d/2,$$

contradicting the choice of d . Thus, $h(\mathbf{x})$ reaches its minimum at some point \mathbf{u} in the open set $B(\mathbf{z}, r)$. Fermat's principle requires that the differential

$$dh(\mathbf{x}) = -2[\mathbf{y} - f(\mathbf{x})]^* df(\mathbf{x}) \quad (4.14)$$

vanish at \mathbf{u} . Given the invertibility of $df(\mathbf{u})$, we therefore have $f(\mathbf{u}) = \mathbf{y}$.

Finally replace V by the open set $B(\mathbf{z}, r) \cap f^{-1}\{B[f(\mathbf{z}), d/2]\}$ contained within it. Our arguments have shown that $f(\mathbf{x})$ is one-to-one from V onto $W = B[f(\mathbf{z}), d/2]$. This allows us to define the inverse function $g(\mathbf{x})$ from W onto V and completes the proof. ■

Example 4.6.1 Polar Coordinates

Consider the transformation

$$\begin{aligned} r &= \|\mathbf{x}\| \\ \theta &= \arctan(x_2/x_1) \end{aligned}$$

to polar coordinates in \mathbb{R}^2 . A brief calculation shows that this transformation has differential

$$\begin{pmatrix} \frac{\partial}{\partial x_1} r & \frac{\partial}{\partial x_2} r \\ \frac{\partial}{\partial x_1} \theta & \frac{\partial}{\partial x_2} \theta \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\frac{1}{r} \sin \theta & \frac{1}{r} \cos \theta \end{pmatrix}.$$

Because the determinant of the differential is $1/r$, the transformation is locally invertible wherever $r \neq 0$. Excluding the half axis $\{x_1 \leq 0, x_2 = 0\}$, the polar transformation maps the plane one-to-one and onto the open set $(0, \infty) \times (-\pi, \pi)$. Here the inverse transformation reduces to

$$\begin{aligned} x_1 &= r \cos \theta \\ x_2 &= r \sin \theta \end{aligned}$$

with differential

$$\begin{pmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{pmatrix} = \begin{pmatrix} \cos \theta & \sin \theta \\ -\frac{1}{r} \sin \theta & \frac{1}{r} \cos \theta \end{pmatrix}^{-1}.$$

We now turn to the implicit function theorem. ■

Proposition 4.6.2 *Let $f(\mathbf{x}, \mathbf{y})$ map an open set $S \subset \mathbb{R}^{m+n}$ into \mathbb{R}^m . Suppose that $f(\mathbf{a}, \mathbf{b}) = \mathbf{0}$ and that $f(\mathbf{x}, \mathbf{y})$ is continuously differentiable on a neighborhood of $(\mathbf{a}, \mathbf{b}) \in S$. If we split the differential*

$$df(\mathbf{x}, \mathbf{y}) = [\partial_1 f(\mathbf{x}, \mathbf{y}), \partial_2 f(\mathbf{x}, \mathbf{y})]$$

into an $m \times m$ block and an $m \times n$ block, and if $\partial_1 f(\mathbf{a}, \mathbf{b})$ is invertible, then there exists a neighborhood U of $\mathbf{b} \in \mathbb{R}^n$ and a continuously differentiable function $g(\mathbf{y})$ from U into \mathbb{R}^m such that $f[g(\mathbf{y}), \mathbf{y}] = \mathbf{0}$. Furthermore, $g(\mathbf{y})$ is unique and has differential

$$dg(\mathbf{y}) = -\partial_1 f[g(\mathbf{y}), \mathbf{y}]^{-1} \partial_2 f[g(\mathbf{y}), \mathbf{y}].$$

Proof: The function

$$h(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} f(\mathbf{x}, \mathbf{y}) \\ \mathbf{y} \end{pmatrix}$$

from S to \mathbb{R}^{m+n} is continuously differentiable with differential

$$dh(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \partial_1 f(\mathbf{x}, \mathbf{y}) & \partial_2 f(\mathbf{x}, \mathbf{y}) \\ \mathbf{0} & \mathbf{I}_n \end{pmatrix}.$$

To apply the inverse function theorem to $h(\mathbf{x}, \mathbf{y})$, we must check that $dh(\mathbf{a}, \mathbf{b})$ is invertible. This is straightforward because

$$\begin{pmatrix} \partial_1 f(\mathbf{a}, \mathbf{b}) & \partial_2 f(\mathbf{a}, \mathbf{b}) \\ \mathbf{0} & \mathbf{I}_n \end{pmatrix} \begin{pmatrix} \mathbf{u} \\ \mathbf{v} \end{pmatrix} = \mathbf{0}$$

can only occur if $\mathbf{v} = \mathbf{0}$ and $\partial_1 f(\mathbf{a}, \mathbf{b})\mathbf{u} = \mathbf{0}$. In view of the invertibility of $\partial_1 f(\mathbf{a}, \mathbf{b})$, the second of these equalities entails $\mathbf{u} = \mathbf{0}$.

Given the invertibility of $dh(\mathbf{a}, \mathbf{b})$, the inverse function theorem implies that $h(\mathbf{x}, \mathbf{y})$ possesses a continuously differentiable inverse that maps an open set W containing $(\mathbf{0}, \mathbf{b})$ onto an open set V containing (\mathbf{a}, \mathbf{b}) . The inverse function takes a point (\mathbf{z}, \mathbf{y}) into the point $[k(\mathbf{z}, \mathbf{y}), \mathbf{y}]$. Consider the function $g(\mathbf{y}) = k(\mathbf{0}, \mathbf{y})$ defined for $(\mathbf{0}, \mathbf{y}) \in W$. Being open, W contains a ball of radius r around $(\mathbf{0}, \mathbf{b})$. There is no harm in restricting the domain of $g(\mathbf{y})$ to the ball $U = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{b}\| < r\}$. On this domain, $g(\mathbf{y})$ is continuously differentiable and $f[g(\mathbf{y}), \mathbf{y}] = \mathbf{0}$. Because $h(\mathbf{x}, \mathbf{y})$ is one-to-one, $g(\mathbf{y})$ is uniquely determined.

Finally, we can identify the differential of $g(\mathbf{y})$ by constructing a slope function. If we let $f(\mathbf{x}, \mathbf{y})$ have slope function $[s_1(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y}), s_2(\mathbf{u}, \mathbf{v}, \mathbf{x}, \mathbf{y})]$ at (\mathbf{x}, \mathbf{y}) , then

$$\begin{aligned} \mathbf{0} &= f[g(\mathbf{v}), \mathbf{v}] - f[g(\mathbf{y}), \mathbf{y}] \\ &= s_1[g(\mathbf{v}), \mathbf{v}, g(\mathbf{y}), \mathbf{y}][g(\mathbf{v}) - g(\mathbf{y})] + s_2[g(\mathbf{v}), \mathbf{v}, g(\mathbf{y}), \mathbf{y}](\mathbf{v} - \mathbf{y}). \end{aligned}$$

The invertibility of $s_1[g(\mathbf{v}), \mathbf{v}, g(\mathbf{y}), \mathbf{y}]$ for \mathbf{v} near \mathbf{y} therefore gives

$$g(\mathbf{v}) - g(\mathbf{y}) = -s_1[g(\mathbf{v}), \mathbf{v}, g(\mathbf{y}), \mathbf{y}]^{-1} s_2[g(\mathbf{v}), \mathbf{v}, g(\mathbf{y}), \mathbf{y}](\mathbf{v} - \mathbf{y}).$$

In the limit as \mathbf{v} approaches \mathbf{y} , we recover the differential $dg(\mathbf{y})$. ■

Example 4.6.2 *Circles and the Implicit Function Theorem*

The function $f(x, y) = x^2 + y^2 - r^2$ has differential $df(x, y) = (2x, 2y)$. Unless $a = 0$, the conditions of the implicit function theorem apply at (a, b) . The choice

$$g(y) = \operatorname{sgn}(a)\sqrt{r^2 - y^2} \quad \text{with} \quad g'(y) = -\frac{y}{x}.$$

clearly satisfies $g(b) = a$ and $g(y)^2 + y^2 - r^2 = 0$. ■

Example 4.6.3 *Tangent Vectors and Tangent Curves*

The intuitive discussion of tangent vectors in Sect. 1.4 can be made rigorous by introducing the notion of a tangent curve at the point \mathbf{x} . This is simply a differentiable curve $v(s)$ having a neighborhood of the scalar 0 as its domain and satisfying $v(0) = \mathbf{x}$ and $g_i[v(s)] = 0$ for all equality constraints and all s sufficiently close to 0. If we apply the chain rule to the composite function $g_i[v(s)] = 0$, then the identity $dg_i(\mathbf{x})v'(0) = 0$ emerges. The vector $\mathbf{w} = v'(0)$ is said to be a tangent vector at \mathbf{x} . Conversely, if \mathbf{w} satisfies $dg_i(\mathbf{x})\mathbf{w} = 0$ for all i , then we can construct a tangent curve at \mathbf{x} with tangent vector \mathbf{w} . This application of the implicit function theorem requires a little notation. Let $G(\mathbf{x})$ be the vector-valued function with i th component $g_i(\mathbf{x})$. The differential $dG(\mathbf{x})$ is the Jacobi matrix whose i th row is the differential $dg_i(\mathbf{x})$. In agreement with our earlier notation, $\nabla G(\mathbf{x})$ is the transpose of $dG(\mathbf{x})$.

Now consider the relationship

$$h(\mathbf{u}, s) = G[\mathbf{x} + \nabla G(\mathbf{x})\mathbf{u} + s\mathbf{w}] = \mathbf{0}.$$

Applying the chain rule to the function $h(\mathbf{u}, s)$ gives

$$\begin{aligned} \partial_{\mathbf{u}}h(\mathbf{0}, 0) &= \left. dG[\mathbf{x} + \nabla G(\mathbf{x})\mathbf{u} + s\mathbf{w}]\nabla G(\mathbf{x}) \right|_{(\mathbf{u}, s)=(\mathbf{0}, 0)} \\ &= dG(\mathbf{x})\nabla G(\mathbf{x}). \end{aligned}$$

Since $G(\mathbf{x}) = \mathbf{0}$ and $dG(\mathbf{x})\nabla G(\mathbf{x})$ is invertible when $dG(\mathbf{x})$ has full row rank, the implicit function theorem implies that we can solve for \mathbf{u} as a function of s in a neighborhood of 0. If we denote the resulting continuously differentiable function by $u(s)$, then our tangent curve is

$$v(s) = \mathbf{x} + \nabla G(\mathbf{x})u(s) + s\mathbf{w}.$$

By definition $u(0) = \mathbf{0}$, $v(0) = \mathbf{x}$, and $G[v(s)] = \mathbf{0}$ for all s close to 0. Thus, we need only check that $v'(0) = \mathbf{w}$. Because

$$v'(0) = \nabla G(\mathbf{x})u'(0) + \mathbf{w},$$

it suffices to check that $\nabla G(\mathbf{x})u'(0) = \mathbf{0}$. However, in view of the equality

$$0 = u'(0)^*\mathbf{0} = u'(0)^*\frac{d}{ds}h[u(0), 0] = u'(0)^*dG(\mathbf{x})[\nabla G(\mathbf{x})u'(0) + \mathbf{w}]$$

and the assumption $dG(\mathbf{x})\mathbf{w} = \mathbf{0}$, this fact is obvious. ■

4.7 Differentials of Matrix-Valued Functions

To define the differential of a matrix-valued function [184], it helps to unroll Carathéodory's definition of differentiability for a vector-valued function $f(\mathbf{x})$. Let us rewrite the slope expansion $f(\mathbf{y}) - f(\mathbf{x}) = s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$ as

$$f(\mathbf{y}) - f(\mathbf{x}) = \sum_{j=1}^m s_j(\mathbf{y}, \mathbf{x})(y_j - x_j) \quad (4.15)$$

using the columns $s_j(\mathbf{y}, \mathbf{x})$ of the slope matrix $s(\mathbf{y}, \mathbf{x})$. This notational retreat retains the linear dependence of the difference $f(\mathbf{y}) - f(\mathbf{x})$ on the increment $\mathbf{y} - \mathbf{x}$ and suggests how to deal with matrix-valued functions. The key step is simply to re-interpret equation (4.15) by replacing the vector-valued function $f(\mathbf{x})$ by a matrix-valued function $f(\mathbf{x})$ and the vector-valued slope $s_j(\mathbf{y}, \mathbf{x})$ by a matrix-valued slope $s_j(\mathbf{y}, \mathbf{x})$. We retain the requirement that $\lim_{\mathbf{y} \rightarrow \mathbf{x}} s_j(\mathbf{y}, \mathbf{x}) = s_j(\mathbf{x}, \mathbf{x})$ for each j . The partial differential matrices $s_j(\mathbf{x}, \mathbf{x}) = \partial_j f(\mathbf{x})$ collectively constitute the differential of $f(\mathbf{x})$. The gratifying thing about this revised definition of differentiability is that it applies to scalars, vectors, and matrices in a unified way. Furthermore, the components of the differential match the scalar, vector, or matrix nature of the original function. We now illustrate the virtue of this perspective by several examples involving matrix differentials.

Example 4.7.1 *The Sum and Transpose Rules*

The rules

$$d[f(\mathbf{x}) + g(\mathbf{x})] = df(\mathbf{x}) + dg(\mathbf{x}), \quad df(\mathbf{x})^* = [df(\mathbf{x})]^*$$

flow directly from the above definition of a differential. ■

Example 4.7.2 *The Chain Rule*

Consider the composition $g \circ f(\mathbf{x})$ of a matrix-valued function with a vector-valued function. Let $g(\mathbf{y})$ have the slope expansion

$$g(\mathbf{y}) - g(\mathbf{x}) = \sum_j t_j(\mathbf{y}, \mathbf{x})(y_j - x_j),$$

and let the j th component $f_j(\mathbf{v})$ of $f(\mathbf{v})$ have the slope expansion

$$f_j(\mathbf{v}) - f_j(\mathbf{u}) = \sum_k s_{jk}(\mathbf{v}, \mathbf{u})(v_k - u_k).$$

Then the identity

$$\begin{aligned} g \circ f(\mathbf{v}) - g \circ f(\mathbf{u}) &= \sum_j t_j[f(\mathbf{v}), f(\mathbf{u})][f_j(\mathbf{v}) - f_j(\mathbf{u})] \\ &= \sum_j t_j[f(\mathbf{v}), f(\mathbf{u})] \sum_k s_{jk}(\mathbf{v}, \mathbf{u})(v_k - u_k) \end{aligned}$$

shows that $g \circ f(\mathbf{v})$ has a differential with $\sum_j \partial_j g[f(\mathbf{u})] \partial_k f_j(\mathbf{u})$ as its k th component. ■

Example 4.7.3 *The Product Rule*

The matrix product $f(\mathbf{y})g(\mathbf{y})$ satisfies the identities

$$\begin{aligned} f(\mathbf{y})g(\mathbf{y}) - f(\mathbf{x})g(\mathbf{x}) &= [f(\mathbf{y}) - f(\mathbf{x})]g(\mathbf{y}) + f(\mathbf{x})[g(\mathbf{y}) - g(\mathbf{x})] \\ &= f(\mathbf{y})[g(\mathbf{y}) - g(\mathbf{x})] + [f(\mathbf{y}) - f(\mathbf{x})]g(\mathbf{x}). \end{aligned}$$

Substituting the slope formulas

$$\begin{aligned} f(\mathbf{y}) - f(\mathbf{x}) &= \sum_k s_k(\mathbf{y}, \mathbf{x})(y_k - x_k) \\ g(\mathbf{y}) - g(\mathbf{x}) &= \sum_k t_k(\mathbf{y}, \mathbf{x})(y_k - x_k), \end{aligned}$$

in either identity demonstrates that $f(\mathbf{y})g(\mathbf{y})$ possesses a differential with k th component $\partial_k f(\mathbf{x})g(\mathbf{x}) + f(\mathbf{x})\partial_k g(\mathbf{x})$ at \mathbf{x} .

When the product $f(\mathbf{x})g(\mathbf{x})$ is a square matrix, we can apply the trace operator to the difference $f(\mathbf{y})g(\mathbf{y}) - f(\mathbf{x})g(\mathbf{x})$ and conclude that

$$\partial_k \operatorname{tr}[f(\mathbf{x})g(\mathbf{x})] = \operatorname{tr}[\partial_k f(\mathbf{x})g(\mathbf{x}) + f(\mathbf{x})\partial_k g(\mathbf{x})].$$

If $f(\mathbf{x})$ is a square matrix, then we can differentiate polynomials in the argument $f(\mathbf{x})$. In view of the sum rule, it suffices to show how to differentiate powers of $f(\mathbf{x})$. For example,

$$\begin{aligned} \partial_k [f(\mathbf{x})]^3 &= \partial_k f(\mathbf{x})[f(\mathbf{x})]^2 + f(\mathbf{x})\partial_k [f(\mathbf{x})]^2 \\ &= \partial_k f(\mathbf{x})[f(\mathbf{x})]^2 + f(\mathbf{x})\partial_k f(\mathbf{x})f(\mathbf{x}) + [f(\mathbf{x})]^2\partial_k f(\mathbf{x}). \end{aligned}$$

In general,

$$\partial_k [f(\mathbf{x})]^n = \partial_k f(\mathbf{x})f(\mathbf{x})^{n-1} + \dots + f(\mathbf{x})^{n-1}\partial_k f(\mathbf{x}).$$

Because matrix multiplication is not commutative, further simplification is impossible. Differentiation of polynomials constitutes a special case of differentiation of power series as discussed in a moment. ■

Example 4.7.4 *The Inverse Rule*

Suppose $f(\mathbf{y})$ is an invertible matrix with slope expansion

$$f(\mathbf{y}) - f(\mathbf{x}) = \sum_k s_k(\mathbf{y}, \mathbf{x})(y_k - x_k).$$

Either of the identities

$$\begin{aligned} f(\mathbf{y})^{-1} - f(\mathbf{x})^{-1} &= -f(\mathbf{x})^{-1}[f(\mathbf{y}) - f(\mathbf{x})]f(\mathbf{y})^{-1} \\ &= f(\mathbf{x})^{-1} \sum_k s_k(\mathbf{y}, \mathbf{x})(y_k - x_k)f(\mathbf{y})^{-1} \end{aligned}$$

$$\begin{aligned} f(\mathbf{y})^{-1} - f(\mathbf{x})^{-1} &= -f(\mathbf{y})^{-1}[f(\mathbf{y}) - f(\mathbf{x})]f(\mathbf{x})^{-1} \\ &= f(\mathbf{y})^{-1} \sum_k s_k(\mathbf{y}, \mathbf{x})(y_k - x_k)f(\mathbf{x})^{-1} \end{aligned}$$

entail a differential of $f(\mathbf{y})^{-1}$ with k th component $-f(\mathbf{x})^{-1}\partial_k f(\mathbf{x})f(\mathbf{x})^{-1}$ at the point \mathbf{x} . \blacksquare

Example 4.7.5 *Matrix Power Series*

Let $p(x) = \sum_{n=0}^{\infty} c_n x^n$ denote a power series with radius of convergence r . As we have seen for the choices $p(x) = e^x$, $p(x) = (1-x)^{-1}$, and $p(x) = \ln(1-x)$, one defines the matrix power series $p(\mathbf{M})$ by substituting a square matrix \mathbf{M} with norm $\|\mathbf{M}\| < r$ for the scalar argument x [128]. On any disc $\{x : |x| \leq s < r\}$ strictly inside its circle of convergence, a power series $p(x)$ is analytic and converges uniformly and absolutely [134, 223]. Furthermore, $p(x)$ can be differentiated term by term; its derivative $p'(x)$ retains r as its radius of convergence.

An obvious question of interest is whether $p(\mathbf{M})$ is differentiable. The easiest route to an affirmative answer exploits forward directional derivatives. Appendix A.6 introduces the notion of a semidifferentiable function. In the current context, this is a function $f(\mathbf{M})$ possessing all possible forward directional derivatives in the strong sense of Hadamard. Thus, we demand the existence of the uniform limit

$$\lim_{t \downarrow 0, \mathbf{U} \rightarrow \mathbf{V}} \frac{f(\mathbf{M} + t\mathbf{U}) - f(\mathbf{M})}{t} = d_{\mathbf{V}}f(\mathbf{M})$$

for all \mathbf{V} . For a semidifferentiable function, Proposition A.6.3 proves that the ordinary differential $df(\mathbf{M})$ exists whenever the map $\mathbf{V} \mapsto d_{\mathbf{V}}f(\mathbf{M})$ is linear. Our experience with polynomials suggests that

$$d_{\mathbf{V}}p(\mathbf{M}) = \sum_{n=1}^{\infty} c_n (\mathbf{V}\mathbf{M}^{n-1} + \dots + \mathbf{M}^{n-1}\mathbf{V}).$$

This series certainly qualifies as linear in \mathbf{V} . It converges because

$$\|\mathbf{V}\mathbf{M}^{n-1} + \dots + \mathbf{M}^{n-1}\mathbf{V}\| \leq n\|\mathbf{V}\| \cdot \|\mathbf{M}\|^{n-1}$$

and the series $\sum_{n=1}^{\infty} n c_n x^{n-1}$ for $p'(x)$ converges absolutely.

Consider the difference quotient

$$\frac{p(\mathbf{M} + t\mathbf{U}) - p(\mathbf{M})}{t} = \sum_{n=1}^{\infty} c_n (\mathbf{U}\mathbf{M}^{n-1} + \dots + \mathbf{M}^{n-1}\mathbf{U}) + \frac{1}{t}R(\mathbf{U}).$$

The series displayed on the right of this equality converges to $d_{\mathbf{V}}p(\mathbf{M})$. Furthermore, the norm of the remainder is bounded above by

$$\begin{aligned}
 \frac{1}{t} \|R(\mathbf{U})\| &\leq \frac{1}{t} \sum_{n=2}^{\infty} c_n \sum_{j=2}^n \binom{n}{j} \|\mathbf{M}\|^{n-j} \|t\mathbf{U}\|^j \\
 &= t \|\mathbf{U}\|^2 \sum_{n=2}^{\infty} c_n \sum_{j=2}^n \frac{n(n-1)}{j(j-1)} \binom{n-2}{j-2} \|\mathbf{M}\|^{n-j} \|t\mathbf{U}\|^{j-2} \\
 &\leq t \|\mathbf{U}\|^2 \sum_{n=2}^{\infty} n(n-1) c_n (\|\mathbf{M}\| + t\|\mathbf{U}\|)^{n-2}.
 \end{aligned}$$

Comparison with the absolutely convergent series $\sum_{n=2}^{\infty} n(n-1)c_n x^{n-2}$ for $p''(x)$ shows that the remainder tends uniformly in norm to $\mathbf{0}$. ■

Example 4.7.6 *Differential of a Determinant*

Sometimes it is simpler to calculate the old-fashioned way with partial derivatives. For example, consider $\det M(\mathbf{x})$. In this case, we exploit the determinant expansion

$$\det \mathbf{M} = \sum_j m_{ij} M_{ij} \tag{4.16}$$

of a square matrix \mathbf{M} in terms of the entries m_{ij} and corresponding cofactors M_{ij} of its i th row. If we ignore the dependence of $M(\mathbf{x})$ on \mathbf{x} and view \mathbf{M} exclusively as a function of its entries m_{ij} , then the expansion (4.16) gives

$$\frac{\partial}{\partial m_{ij}} \det \mathbf{M} = M_{ij}.$$

The chain rule therefore implies

$$\begin{aligned}
 \frac{\partial}{\partial x_i} \det M(\mathbf{x}) &= \sum_j \sum_k \frac{\partial}{\partial m_{jk}} \det M(\mathbf{x}) \frac{\partial}{\partial x_i} m_{jk}(\mathbf{x}) \\
 &= \sum_j \sum_k M_{jk}(\mathbf{x}) \frac{\partial}{\partial x_i} m_{jk}(\mathbf{x}).
 \end{aligned}$$

According to Cramer’s rule, this can be simplified by noting that the matrix with entry $(\det \mathbf{M})^{-1} M_{jk}$ in row k and column j is \mathbf{M}^{-1} . It follows that

$$\frac{\partial}{\partial x_i} \det M(\mathbf{x}) = \det M(\mathbf{x}) \operatorname{tr} \left[M(\mathbf{x})^{-1} \frac{\partial}{\partial x_i} M(\mathbf{x}) \right]$$

when $M(\mathbf{x})$ is invertible. If $\det M(\mathbf{x})$ is positive, for instance if $M(\mathbf{x})$ is positive definite, then we have the even cleaner formula

$$\frac{\partial}{\partial x_i} \ln \det M(\mathbf{x}) = \operatorname{tr} \left[M(\mathbf{x})^{-1} \frac{\partial}{\partial x_i} M(\mathbf{x}) \right].$$

This is sometimes expressed in the ambiguous but suggestive form

$$d[\ln \det M(\mathbf{x})] = \operatorname{tr} [M(\mathbf{x})^{-1} dM(\mathbf{x})]. \quad (4.17)$$

As an application of these results, consider minimization of the function

$$f(\mathbf{M}) = a \ln \det \mathbf{M} + b \operatorname{tr}(\mathbf{S}\mathbf{M}^{-1}), \quad (4.18)$$

where a and b are positive constants and \mathbf{S} and \mathbf{M} are positive definite matrices. On its open domain $f(\mathbf{M})$ is obviously differentiable. Because \mathbf{M} is symmetric, we parameterize it by its lower triangle, including of course its diagonal. At a local minimum, the differential of $f(\mathbf{M})$ vanishes. In view of the cyclic permutation property of the trace function, this differential has components

$$\begin{aligned} \partial_k f(\mathbf{M}) &= a \operatorname{tr}(\mathbf{M}^{-1} \partial_k \mathbf{M}) - b \operatorname{tr}(\mathbf{S}\mathbf{M}^{-1} \partial_k \mathbf{M}\mathbf{M}^{-1}) \\ &= \operatorname{tr}[\partial_k \mathbf{M}(a\mathbf{M}^{-1} - b\mathbf{M}^{-1}\mathbf{S}\mathbf{M}^{-1})] \\ &= \operatorname{tr}(\partial_k \mathbf{M}\mathbf{N}) \end{aligned}$$

in obvious notation. If k corresponds to the i th diagonal entry of \mathbf{M} , then $\partial_k \mathbf{M}$ has all entries 0 except for a 1 in the i th diagonal entry. If k corresponds to an off-diagonal entry in row i and column j , then $\partial_k \mathbf{M}$ has all entries 0 except for a 1 in the off-diagonal entries in the symmetric positions (i, j) and (j, i) . An easy calculation based on $\partial_k f(\mathbf{M}) = \operatorname{tr}(\partial_k \mathbf{M}\mathbf{N})$ now shows that

$$\partial_k f(\mathbf{M}) = \begin{cases} n_{ii} & k = (i, i) \\ 2n_{ij} & k = (i, j), \quad i > j. \end{cases}$$

Because the components of the differential vanish, the matrix \mathbf{N} must vanish as well. In other words $\mathbf{N} = a\mathbf{M}^{-1} - b\mathbf{M}^{-1}\mathbf{S}\mathbf{M}^{-1} = \mathbf{0}$. The one and only solution to this equation is $\mathbf{M} = \frac{b}{a}\mathbf{S}$. In Example 6.3.12 we will prove that this stationary point provides the minimum of $f(\mathbf{M})$. ■

4.8 Problems

1. Verify the entries in Table 4.1 not derived in the text.
2. For each positive integer n and real number x , find the derivative, if possible, of the function

$$f_n(x) = \begin{cases} x^n \sin\left(\frac{1}{x}\right) & x \neq 0 \\ 0 & x = 0. \end{cases}$$

Pay particular attention to the point 0.

3. Show that the function

$$f(x) = x \ln(1 + x^{-1})$$

is strictly increasing on $(0, \infty)$ and satisfies $\lim_{x \rightarrow 0} f(x) = 0$ and $\lim_{x \rightarrow \infty} f(x) = 1$ [69].

4. Let $h(x) = f(x)g(x)$ be the product of two functions that are each k times differentiable. Derive Leibnitz's formula

$$h^{(k)}(x) = \sum_{j=0}^k \binom{k}{j} f^{(j)}(x)g^{(k-j)}(x).$$

5. Assume that the real-valued functions $f(y)$ and $g(y)$ are differentiable at the real point x . If (a) $f(x) = g(x) = 0$, (b) $g(y) \neq 0$ for y near x , and (c) $g'(x) \neq 0$, then demonstrate L'Hôpital's rule

$$\lim_{y \rightarrow x} \frac{f(y)}{g(y)} = \frac{f'(x)}{g'(x)}.$$

6. Let $f(x)$ and $g(x)$ be continuous on the closed interval $[a, b]$ and differentiable on the open interval (a, b) . Prove that there exists a point $x \in (a, b)$ such that

$$[f(b) - f(a)]g'(x) = [g(b) - g(a)]f'(x).$$

(Hint: Consider $h(x) = [f(b) - f(a)]g(x) - [g(b) - g(a)]f(x)$.)

7. Show that $\frac{\sin x}{x} < 1$ for $x \neq 0$. Use this fact to prove $1 - \frac{x^2}{2} < \cos x$ for $x \neq 0$.
8. Suppose the positive function $f(x)$ satisfies the inequality

$$f'(x) \leq cf(x)$$

for some constant $c \geq 0$ and all $x \geq 0$. Prove that $f(x) \leq e^{cx} f(0)$ for $x \geq 0$ [69].

9. Abbreviate the scalar exponential and logarithmic functions by $e(t)$ and $l(t)$. Based on the defining differential equations $e'(t) = e(t)$ and $l'(t) = t^{-1}$ and the initial conditions $e(0) = 1$ and $l(1) = 0$, prove that $e \circ l(t) = t$ for $t > 0$ and $l \circ e(t) = t$ for all t . Thus, $e(t)$ and $l(t)$ are functional inverses. (Hint: The derivatives of $l \circ e(t)$ and $t^{-1}e \circ l(t)$ are constant.)

10. Prove the identities $\cos x = \cos(-x)$ and $\sin x = -\sin(-x)$ and the identities

$$\begin{aligned}\cos(x+y) &= \cos x \cos y - \sin x \sin y \\ \sin(x+y) &= \sin x \cos y + \cos x \sin y.\end{aligned}$$

(Hint: The differential equation

$$f'(x) = \begin{pmatrix} 0 & -1 \\ 1 & 0 \end{pmatrix} f(x)$$

with $f(0)$ fixed has a unique solution. In each case demonstrate that both sides of the two proposed identities satisfy the differential equation. The initial condition is $f(0) = (1, 0)^*$ in the first case and $f(0) = (\cos y, \sin y)^*$ in the second case.)

11. Use the defining differential equations for $\cos x$ and $\sin x$ to show that there is a smallest positive root $\frac{\pi}{2}$ of the equation $\cos x = 0$. Applying the trigonometric identities of the previous problem, deduce that $\cos \pi = -1$ and that $\cos(x+2\pi) = \cos x$ and $\sin(x+2\pi) = \sin x$. (Hint: Argue by contradiction that $\cos x > 0$ cannot hold for all positive x .)
12. Let $f(x, y)$ be an integrable function of x for each y . If the partial derivative $\frac{\partial}{\partial y} f(x, y)$ exists, it makes sense to ask when the interchange

$$\frac{d}{dy} \int_a^b f(x, y) dx = \int_a^b \frac{\partial}{\partial y} f(x, y) dx$$

is permissible. Demonstrate that a sufficient condition is the existence of integrable functions $g(x)$ and $h(x)$ satisfying

$$g(x) \leq \frac{\partial}{\partial y} f(x, y) \leq h(x)$$

for all x and y . Show how to construct $g(x)$ and $h(x)$ when $\frac{\partial}{\partial y} f(x, y)$ is jointly continuous in x and y . (Hint: Apply the mean value theorem to the difference quotient. Then invoke the dominated convergence theorem.)

13. Let $f(x)$ be a differentiable curve mapping the interval $[a, b]$ into \mathbb{R}^n . Show that $\|f(x)\|$ is constant if and only if $f(x)$ and $f'(x)$ are orthogonal for all x .
14. Suppose the constants c_0, \dots, c_k satisfy the condition

$$c_0 + \frac{c_1}{2} + \dots + \frac{c_k}{k+1} = 0.$$

Demonstrate that the polynomial $p(x) = c_0 + c_1 x + \dots + c_k x^k$ has a root on the interval $[0, 1]$.

15. Prove that the function

$$f(x) = \begin{cases} e^{-x^{-2}} & x \neq 0 \\ 0 & x = 0 \end{cases}$$

is infinitely differentiable and has derivative $f^{(n)}(0) = 0$ for every n .

16. Consider the ordinary differential equation $M'(t) = N(t)M(t)$ with initial condition $M(0) = \mathbf{A}$ for $n \times n$ matrices. If \mathbf{A} is invertible, then demonstrate that any two solutions coincide in a neighborhood of 0. (Hint: If $P(t)$ and $Q(t)$ are two solutions, then differentiate the product $P(t)^{-1}Q(t)$ using the product and inverse rules.)
17. Let $f(\mathbf{x})$ and $g(\mathbf{x})$ be real-valued functions defined on a neighborhood of \mathbf{y} in \mathbb{R}^n . If $f(\mathbf{x})$ is differentiable at \mathbf{y} and has $f(\mathbf{y}) = 0$, and $g(\mathbf{x})$ is continuous at \mathbf{y} , then prove that $f(\mathbf{x})g(\mathbf{x})$ is differentiable at \mathbf{y} .
18. Consider a continuous function $f(x, y)$ defined on a compact interval $I = [a, b] \times [c, d]$ of \mathbb{R}^2 . Assume that the partial derivative $\partial_2 f(x, y)$ is continuous on I and that $p(y)$ and $q(y)$ are differentiable functions mapping $[c, d]$ into $[a, b]$. In this setting, prove that the integral

$$F(y) = \int_{p(y)}^{q(y)} f(x, y) dx$$

has derivative

$$F'(y) = \int_{p(y)}^{q(y)} \partial_2 f(x, y) dx + f[q(y), y]q'(y) - f[p(y), y]p'(y).$$

(Hint: Use Problem 12.)

19. A real-valued function $f(\mathbf{x})$ on \mathbb{R}^n is said to be homogeneous of integer order $k \geq 1$ if $f(t\mathbf{x}) = t^k f(\mathbf{x})$ for every scalar t and $\mathbf{x} \in \mathbb{R}^n$. For instance, $f(x) = x^3$ is homogeneous of order 3 on \mathbb{R} . Demonstrate that a homogeneous differentiable function $f(\mathbf{x})$ satisfies $df(\mathbf{x})\mathbf{x} = kf(\mathbf{x})$.
20. As a converse to the chain rule stated in Proposition 4.4.3, suppose (a) $f(\mathbf{y})$ is continuous at \mathbf{x} , (b) $g(\mathbf{z})$ is differentiable at $f(\mathbf{x})$, (c) $g \circ f(\mathbf{y})$ is differentiable at \mathbf{x} , and (d) the matrix $dg[f(\mathbf{x})]$ is invertible. Show that $f(\mathbf{y})$ is differentiable at \mathbf{x} . (Hint: Equate the two slope expansions of $g \circ f(\mathbf{y}) - g \circ f(\mathbf{x})$.)
21. A real-valued function $f(\mathbf{y})$ is said to be Gâteaux differentiable at \mathbf{x} if there exists a vector \mathbf{g} such that

$$\lim_{t \rightarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \mathbf{g}^* \mathbf{v}$$

for all unit vectors \mathbf{v} . In other words, all directional derivatives exist and depend linearly on the direction \mathbf{v} . In general, Gâteaux differentiability does not imply differentiability. However, if $f(\mathbf{y})$ satisfies a Lipschitz condition $|f(\mathbf{y}) - f(\mathbf{z})| \leq c\|\mathbf{y} - \mathbf{z}\|$ in a neighborhood of \mathbf{x} , then prove that Gâteaux differentiability at \mathbf{x} implies differentiability at \mathbf{x} with $\nabla f(\mathbf{x}) = \mathbf{g}$. In Chap. 6, the proof of Proposition 6.4.1 shows that a convex function is locally Lipschitz around each of its interior points. Thus, Gâteaux differentiability and differentiability are equivalent for a convex function at an interior point of its domain. Consult Appendix A.6 for a fuller treatment of this topic. (Hint: Every \mathbf{u} on the unit sphere of \mathbb{R}^n is within $\epsilon > 0$ of some member of a finite set $\{\mathbf{u}_1, \dots, \mathbf{u}_m\}$ of points on the sphere. Now write

$$\begin{aligned} f(\mathbf{x} + \mathbf{w}) - f(\mathbf{x}) - \mathbf{g}^* \mathbf{w} &= f(\mathbf{x} + \|\mathbf{w}\| \mathbf{u}_k) - f(\mathbf{x}) - \|\mathbf{w}\| \mathbf{g}^* \mathbf{u}_k \\ &\quad + f(\mathbf{x} + \mathbf{w}) - f(\mathbf{x} + \|\mathbf{w}\| \mathbf{u}_k) \\ &\quad + \|\mathbf{w}\| \mathbf{g}^* \left(\mathbf{u}_k - \frac{\mathbf{w}}{\|\mathbf{w}\|} \right) \end{aligned}$$

for an appropriate choice of \mathbf{u}_k .)

22. Continuing Problem 21, show that the function

$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } x_1 = x_2^2 \text{ and } x_2 \neq 0 \\ 0 & \text{otherwise} \end{cases}$$

is Gâteaux differentiable at the origin $\mathbf{0}$ of \mathbb{R}^2 but not differentiable or even continuous there.

23. Demonstrate that the function

$$f(\mathbf{x}) = \begin{cases} \frac{x_1 x_2}{x_1^2 + x_2} & x_1^2 + x_2 \neq 0 \\ 0 & x_1^2 + x_2 = 0 \end{cases}$$

is discontinuous at $\mathbf{0}$, yet the directional derivative $d_{\mathbf{v}} f(\mathbf{0})$ exists for all \mathbf{v} . In this instance, show that $\mathbf{v} \mapsto d_{\mathbf{v}} f(\mathbf{0})$ is nonlinear in \mathbf{v} .

24. On the set $\{\mathbf{x} \neq \mathbf{0}\}$, demonstrate that $\|\mathbf{x}\|$ is differentiable with differential $d\|\mathbf{x}\| = \|\mathbf{x}\|^{-1} \mathbf{x}^*$.
25. Prove that the directional derivative $d_{\mathbf{v}} \|\mathbf{x}\|$ fails to exist when $\mathbf{x} = \mathbf{0}$ and $\mathbf{v} \neq \mathbf{0}$. Note that the forward directional derivative

$$\lim_{t \downarrow 0} \frac{\|t\mathbf{v}\| - 0}{t} = \|\mathbf{v}\|$$

does exist for all \mathbf{v} .

26. Suppose the vector-valued function $\mathbf{x}(t)$ is differentiable in the scalar t . Show that $\|\mathbf{x}(t)\|' \leq \|\mathbf{x}'(t)\|$ whenever $\mathbf{x}(t) \neq \mathbf{0}$.

27. Continuing Problem 2, show that $f_4(x)$ is continuously differentiable on \mathbb{R} , that its image is an open subset of \mathbb{R} , and yet $f_4(x)$ is not one-to-one on any interval around 0. What bearing does this have on the inverse function theorem?

28. Demonstrate that the function

$$f(x) = \begin{cases} x + 2x^2 \sin\left(\frac{1}{x}\right) & x \neq 0 \\ 0 & x = 0 \end{cases}$$

has a bounded derivative $f'(x)$ for all x and $f'(0) \neq 0$, yet $f(x)$ is not one-to-one on any interval around 0. Why does the inverse function theorem fail in this case?

29. Consider the equation $f(x) = tg(x)$ determined by the continuously differentiable functions $f(x)$ and $g(x)$ from \mathbb{R} into \mathbb{R} . If $f(0) = 0$ and $f'(0) \neq 0$, then show that in a suitably small interval $|t| < \delta$ there is a unique continuously differentiable function $x(t)$ solving the equation and satisfying $x(0) = 0$. Prove that $x'(0) = g(0)/f'(0)$.

30. Suppose that the differential $df(\mathbf{x})$ of the continuously differentiable function $f(\mathbf{x})$ from \mathbb{R}^m to \mathbb{R}^n has full rank at \mathbf{y} . Show that $f(\mathbf{x})$ is one-to-one in a neighborhood of \mathbf{y} . Note that $m \leq n$ must hold.

31. The mean value inequality (4.12) can be improved. Suppose that along the line segment $\{\mathbf{u} = \mathbf{x} + t(\mathbf{y} - \mathbf{x}) : t \in [0, 1]\}$ the gradient $\nabla f(\mathbf{u})$ satisfies the Lipschitz inequality

$$\|\nabla f(\mathbf{u}) - \nabla f(\mathbf{v})\| \leq \lambda \|\mathbf{u} - \mathbf{v}\| \quad (4.19)$$

for some constant $\lambda \geq 0$. This is the case if the second differential $d^2f(\mathbf{u})$ exists and is continuous in \mathbf{u} . Prove that

$$\|f(\mathbf{y}) - f(\mathbf{x}) - df(\mathbf{x})(\mathbf{y} - \mathbf{x})\| \leq \frac{\lambda}{2} \|\mathbf{y} - \mathbf{x}\|^2.$$

(Hint: Invoke the fundamental theorem of calculus, and bound the norm of the integrand by inequality (4.19).)

32. Demonstrate that the bound in Problem 31 can be generalized to

$$\|f(\mathbf{u}) - f(\mathbf{v}) - df(\mathbf{x})(\mathbf{u} - \mathbf{v})\| \leq \frac{\lambda}{2} (\|\mathbf{u} - \mathbf{x}\| + \|\mathbf{v} - \mathbf{x}\|) \|\mathbf{u} - \mathbf{v}\|$$

for any triple of points \mathbf{u} , \mathbf{v} , and \mathbf{x} contained in the convex domain of the function $f(\mathbf{x})$. Note the similarity to the straddle inequality (3.8).

33. Given the result of Problem 32, suppose that the domain and range of $f(\mathbf{x})$ are both contained in \mathbb{R}^n and that the matrix $df(\mathbf{x})$ is invertible. Show that there exist positive constants α , β , and ϵ such that

$$\alpha \|\mathbf{u} - \mathbf{v}\| \leq \|f(\mathbf{u}) - f(\mathbf{v})\| \leq \beta \|\mathbf{u} - \mathbf{v}\|$$

for all \mathbf{u} and \mathbf{v} with $\max\{\|\mathbf{u} - \mathbf{x}\|, \|\mathbf{v} - \mathbf{x}\|\} \leq \epsilon$. (Hints: Write

$$f(\mathbf{u}) - f(\mathbf{v}) = f(\mathbf{u}) - f(\mathbf{v}) - df(\mathbf{x})(\mathbf{u} - \mathbf{v}) + df(\mathbf{x})(\mathbf{u} - \mathbf{v}),$$

and apply the bound $\|\mathbf{u} - \mathbf{v}\| \leq \|df(\mathbf{x})^{-1}\| \cdot \|df(\mathbf{x})(\mathbf{u} - \mathbf{v})\|$.)

34. Consider the function $f(\mathbf{M}) = \mathbf{M} + \mathbf{M}^2$ defined on $n \times n$ matrices. Show that the range of $f(\mathbf{M})$ contains a neighborhood of the trivial matrix $\mathbf{0}$ [69]. (Hint: Compute the differential of $f(\mathbf{M})$ and apply the inverse function theorem.)
35. Calculate the differentials of the matrix-valued functions $e^{\mathbf{M}}$ and $\ln(\mathbf{I} - \mathbf{M})$. What are the values of these differentials at $\mathbf{M} = \mathbf{0}$?
36. For a differentiable matrix $M(\mathbf{x})$, verify the partial derivatives

$$\partial_k(MM^*) = (\partial_k M)M^* + M(\partial_k M)^*$$

$$\partial_k(M^*M) = (\partial_k M)^*M + M^*(\partial_k M)$$

$$\partial_k M^p = \sum_{j=1}^p M^{j-1} \partial_k M M^{p-j}, \quad p > 0$$

$$\partial_k M^{-p} = -\sum_{j=1}^p M^{-j} \partial_k M M^{-p+j-1}, \quad p > 0$$

$$\partial_k \operatorname{tr}(MM^*) = 2 \operatorname{tr}(M^* \partial_k M)$$

$$\partial_k \operatorname{tr}(M^*M) = 2 \operatorname{tr}(M^* \partial_k M)$$

$$\partial_k \operatorname{tr}(M^p) = p \operatorname{tr}(M^{p-1} \partial_k M)$$

$$\partial_k \det(MM^*) = 2 \det(MM^*) \operatorname{tr}[M^*(MM^*)^{-1} \partial_k M]$$

$$\partial_k \det(M^*M) = 2 \det(M^*M) \operatorname{tr}[(M^*M)^{-1} M^* \partial_k M]$$

$$\partial_k \det(M^p) = p \det(M^p) \operatorname{tr}(M^{-1} \partial_k M),$$

where p is an integer and $\partial_k = \frac{\partial}{\partial x_k}$.

37. A random $m \times m$ matrix \mathbf{W} with density

$$f(\mathbf{w}) = \frac{|\det \mathbf{w}|^{(n-m-1)/2} e^{-\frac{1}{2} \operatorname{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{w})}}{2^{nm/2} \pi^{m(m-1)/4} |\det \boldsymbol{\Sigma}|^{n/2} \prod_{i=1}^m \Gamma[(n+1-i)/2]}$$

is said to be Wishart distributed [3]. Here $\boldsymbol{\Sigma}$ is positive definite and n is an integer with $n > m$. It turns out that the inverse matrix $\mathbf{V} = \mathbf{W}^{-1}$ has density

$$g(\mathbf{v}) = \frac{|\det \mathbf{v}|^{-(n+m+1)/2} e^{-\frac{1}{2} \operatorname{tr}(\boldsymbol{\Sigma}^{-1} \mathbf{v}^{-1})}}{2^{nm/2} \pi^{m(m-1)/4} |\det \boldsymbol{\Sigma}|^{n/2} \prod_{i=1}^m \Gamma[(n+1-i)/2]}.$$

This naturally is called the inverse Wishart density. Demonstrate that the modes of the Wishart and inverse Wishart densities occur at $\mathbf{w} = (n - m - 1)\mathbf{\Sigma}$ and $\mathbf{v} = (n + m + 1)^{-1}\mathbf{\Sigma}^{-1}$, respectively. (Hints: Show that these points are stationary points of the log densities by considering the function (4.18).)

5

Karush-Kuhn-Tucker Theory

5.1 Introduction

In the current chapter, we study the problem of minimizing a real-valued function $f(\mathbf{x})$ subject to the constraints

$$\begin{aligned}g_i(\mathbf{x}) &= 0, & 1 \leq i \leq p \\h_j(\mathbf{x}) &\leq 0, & 1 \leq j \leq q.\end{aligned}$$

All of these functions share some open set $U \subset \mathbb{R}^n$ as their domain. Maximizing $f(\mathbf{x})$ is equivalent to minimizing $-f(\mathbf{x})$, so there is no loss of generality in considering minimization. The function $f(\mathbf{x})$ is called the objective function, the functions $g_i(\mathbf{x})$ are called equality constraints, and the functions $h_j(\mathbf{x})$ are called inequality constraints. Any point $\mathbf{x} \in U$ satisfying all of the constraints is said to be feasible. A constraint $h_j(\mathbf{x})$ is active at the feasible point \mathbf{x} provided $h_j(\mathbf{x}) = 0$; it is inactive if $h_j(\mathbf{x}) < 0$. In general, we will assume that the feasible region is nonempty. The case $p = 0$ of no equality constraints and the case $q = 0$ of no inequality constraints are both allowed.

In exploring solutions to the above constrained minimization problem, we will meet a generalization of the Lagrange multiplier rule fashioned independently by Karush [149] and Kuhn and Tucker [161]. Under fairly weak regularity conditions, the rule holds at all extrema. In contrast to this necessary condition, sufficient conditions for an extremum involve second derivatives. To state and prove the most useful sufficient condition, we

must confront second differentials and what it means for a function to be twice differentiable. The matter is straightforward conceptually but computationally messy. Fortunately, we can build on the material presented in Chap. 4.

5.2 The Multiplier Rule

Before embarking on the long and interesting proof of the multiplier rule, we turn to linear programming as a specific example of constrained optimization. A huge literature has grown up around this single application.

Example 5.2.1 Linear Programming

If the objective function $f(\mathbf{x})$ and the constraints $g_i(\mathbf{x})$ and $h_j(\mathbf{x})$ are all affine functions $\mathbf{z}^* \mathbf{x} + c$, then the constrained minimization problem is termed linear programming. In the literature on linear programming, the standard linear program is posed as one of minimizing $f(\mathbf{x}) = \mathbf{z}^* \mathbf{x}$ subject to the linear equality constraints $\mathbf{V} \mathbf{x} = \mathbf{d}$ and the nonnegativity constraints $x_i \geq 0$ for all $1 \leq i \leq n$. The inequality constraints are collectively abbreviated as $\mathbf{x} \geq \mathbf{0}$. To show that the standard linear program encompasses our apparently more general version of linear programming, we note first that we can omit the affine constant in the objective function $f(\mathbf{x})$. The p linear equality constraints

$$0 = g_i(\mathbf{x}) = \mathbf{v}_i^* \mathbf{x} - d_i$$

are already in the form $\mathbf{V} \mathbf{x} = \mathbf{d}$ if we define \mathbf{V} to be the $p \times n$ matrix with i th row \mathbf{v}_i^* and \mathbf{d} to be the $p \times 1$ vector with i th entry d_i . The inequality constraint $h_j(\mathbf{x}) \leq 0$ can be elevated to an equality constraint $h_j(\mathbf{x}) + y_j = 0$ by introducing an additional variable y_j called a slack variable with the stipulation that $y_j \geq 0$. If any of the variables x_i is not already constrained by $x_i \geq 0$, then we can introduce what are termed free variables $u_i \geq 0$ and $w_i \geq 0$ so that $x_i = u_i - w_i$ and replace x_i everywhere by this difference. ■

In proving the multiplier rule, it turns out that one must restrict the behavior of the constraints at a local extremum to avoid redundant constraints. There are several constraint qualifications.

Definition 5.2.1 Mangasarian-Fromovitz Constraint Qualification

This condition holds at a feasible point \mathbf{x} provided the differentials $dg_i(\mathbf{x})$ are linearly independent and there exists a vector \mathbf{v} with $dg_i(\mathbf{x})\mathbf{v} = 0$ for all i and $dh_j(\mathbf{x})\mathbf{v} < 0$ for all inequality constraints $h_j(\mathbf{x})$ active at \mathbf{x} [186]. The vector \mathbf{v} is a tangent vector in the sense that infinitesimal motion from \mathbf{x} along \mathbf{v} stays within the feasible region. ■

Because the Mangasarian-Fromovitz condition is difficult to check, we will consider the simpler sufficient condition of Kuhn and Tucker [161] in the next section. In the meantime, we state and prove the Lagrange multiplier rule extended to inequality constraints by Karush and Kuhn and Tucker. Our proof reproduces McShane's lovely argument, which substitutes penalties for constraints [116, 194].

Proposition 5.2.1 *Suppose the objective function $f(\mathbf{y})$ of the constrained optimization problem has a local minimum at the feasible point \mathbf{x} . If $f(\mathbf{y})$ and the various constraint functions are continuously differentiable near \mathbf{x} , then there exists a unit vector of Lagrange multipliers $\lambda_0, \dots, \lambda_p$ and μ_1, \dots, μ_q such that*

$$\lambda_0 \nabla f(\mathbf{x}) + \sum_{i=1}^p \lambda_i \nabla g_i(\mathbf{x}) + \sum_{j=1}^q \mu_j \nabla h_j(\mathbf{x}) = \mathbf{0}. \quad (5.1)$$

Moreover, each of the multipliers λ_0 and μ_j is nonnegative, and $\mu_j = 0$ whenever $h_j(\mathbf{x}) < 0$. If the constraint functions satisfy the Mangasarian-Fromovitz constraint qualification at \mathbf{x} , then we can take $\lambda_0 = 1$.

Proof: Without loss of generality, we assume $\mathbf{x} = \mathbf{0}$ and $f(\mathbf{0}) = 0$. By renumbering the inequality constraints if necessary, we also suppose that the first r of them are active at $\mathbf{0}$ and the last $q - r$ of them are inactive at $\mathbf{0}$. Now choose $\delta > 0$ so that (a) the closed ball

$$C(\mathbf{0}, \delta) = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y}\| \leq \delta\}$$

is contained in the open domain U , (b) $\mathbf{0}$ is the minimum point of $f(\mathbf{y})$ in $C(\mathbf{0}, \delta)$ subject to the constraints, (c) the objective and constraint functions are continuously differentiable in $C(\mathbf{0}, \delta)$, and (d) the constraints $h_j(\mathbf{x})$ inactive at $\mathbf{0}$ are inactive throughout $C(\mathbf{0}, \delta)$.

On the road to our ultimate goal, consider the functions

$$h_{j+}(\mathbf{y}) = \max\{h_j(\mathbf{y}), 0\}.$$

Using these functions, we now prove that for each $0 < \epsilon \leq \delta$, there exists an $\alpha > 0$ such that

$$f(\mathbf{y}) + \|\mathbf{y}\|^2 + \alpha \sum_{i=1}^p g_i(\mathbf{y})^2 + \alpha \sum_{j=1}^r h_{j+}(\mathbf{y})^2 > 0 \quad (5.2)$$

for all \mathbf{y} with $\|\mathbf{y}\| = \epsilon$. This is not an entirely trivial claim to prove because $f(\mathbf{y})$ can be negative on $C(\mathbf{0}, \delta)$ outside the feasible region.

Suppose the claim is false. Then there is a sequence of points \mathbf{y}_m with $\|\mathbf{y}_m\| = \epsilon$ and a sequence of numbers α_m tending to ∞ such that

$$f(\mathbf{y}_m) + \|\mathbf{y}_m\|^2 \leq -\alpha_m \sum_{i=1}^p g_i(\mathbf{y}_m)^2 - \alpha_m \sum_{j=1}^r h_{j+}(\mathbf{y}_m)^2. \quad (5.3)$$

Because the boundary of $C(\mathbf{0}, \epsilon)$ is compact, the sequence \mathbf{y}_m has a convergent subsequence, which without loss of generality we take to be the sequence itself. The limit \mathbf{z} of the sequence clearly has norm $\|\mathbf{z}\| = \epsilon$. Dividing both sides of inequality (5.3) by $-\alpha_m$ and sending m to ∞ produce

$$\sum_{i=1}^p g_i(\mathbf{z})^2 + \sum_{j=1}^r h_{j+}(\mathbf{z})^2 = 0. \quad (5.4)$$

It follows that \mathbf{z} is feasible with $f(\mathbf{z}) \geq f(\mathbf{0}) = 0$. However, inequality (5.3) requires that each $f(\mathbf{y}_m) \leq -\epsilon^2$. Since this last relation is preserved in the limit, we reach a contradiction and consequently establish the validity of inequality (5.2).

Our next goal is to prove that there exists a point \mathbf{u} and a unit vector $(\lambda_0, \lambda_1, \dots, \lambda_p, \mu_1, \dots, \mu_r)^*$ such that (a) $\|\mathbf{u}\| < \epsilon$, (b) each of the multipliers λ_0 and μ_j is nonnegative, and (c)

$$\lambda_0[\nabla f(\mathbf{u}) + 2\mathbf{u}] + \sum_{i=1}^p \lambda_i \nabla g_i(\mathbf{u}) + \sum_{j=1}^r \mu_j \nabla h_j(\mathbf{u}) = \mathbf{0}. \quad (5.5)$$

Observe here that the distinction between active and inactive constraints comes into play again. To prove the Lagrange multiplier rule (5.5), define

$$F(\mathbf{y}) = f(\mathbf{y}) + \|\mathbf{y}\|^2 + \alpha \sum_{i=1}^p g_i(\mathbf{y})^2 + \alpha \sum_{j=1}^r h_{j+}(\mathbf{y})^2$$

using the α satisfying condition (5.2). For typographical reasons, we omit the dependence of α on ϵ .

Given that $F(\mathbf{y})$ is continuous, there is a point \mathbf{u} giving the unconstrained minimum of $F(\mathbf{y})$ on the compact set $C(\mathbf{0}, \epsilon)$. Because this point satisfies $F(\mathbf{u}) \leq F(\mathbf{0}) = 0$, it is impossible that $\|\mathbf{u}\| = \epsilon$ in view of inequality (5.2). Thus, \mathbf{u} falls in the interior of $C(\mathbf{0}, \epsilon)$ where $\nabla F(\mathbf{u}) = \mathbf{0}$ must occur. The gradient condition $\nabla F(\mathbf{u}) = \mathbf{0}$ can be expressed as

$$\nabla f(\mathbf{u}) + 2\mathbf{u} + \alpha \sum_{i=1}^p 2g_i(\mathbf{u})\nabla g_i(\mathbf{u}) + \alpha \sum_{j=1}^r 2h_{j+}(\mathbf{u})\nabla h_j(\mathbf{u}) = \mathbf{0}, \quad (5.6)$$

invoking the differentiability of the functions $h_{j+}(\mathbf{y})^2$ derived in Example 4.4.3 of Chap. 4. If we divide equality (5.6) by the norm of the vector

$$\mathbf{v} = [1, 2\alpha g_1(\mathbf{u}), \dots, 2\alpha g_p(\mathbf{u}), 2\alpha h_{1+}(\mathbf{u}), \dots, 2\alpha h_{r+}(\mathbf{u})]^*$$

and redefine the Lagrange multipliers accordingly, then the multiplier rule (5.5) holds with each of the multipliers λ_0 and μ_j nonnegative.

Now choose a sequence $\epsilon_m > 0$ tending to 0 and corresponding points \mathbf{u}_m where the Lagrange multiplier rule (5.5) holds. The sequence of unit

vectors $(\lambda_{m0}, \dots, \lambda_{mp}, \mu_{m1}, \dots, \mu_{mr})^*$ has a convergent subsequence with limit $(\lambda_0, \dots, \lambda_p, \mu_1, \dots, \mu_r)^*$ that is also a unit vector. Replacing the sequence \mathbf{u}_m by the corresponding subsequence \mathbf{u}_{m_k} allows us to take limits along \mathbf{u}_{m_k} in equality (5.5) and achieve equality (5.1). Observe that \mathbf{u}_{m_k} converges to $\mathbf{0}$ because $\|\mathbf{u}_{m_k}\| \leq \epsilon_{m_k}$.

Finally, suppose the constraint qualification holds at the local minimum $\mathbf{0}$ and that $\lambda_0 = 0$. If all of the nonnegative multipliers μ_j are 0, then at least one of the λ_i with $1 \leq i \leq p$ is not 0. But this contradicts the linear independence of the $dg_i(\mathbf{0})$. Now consider the vector \mathbf{v} guaranteed by the constraint qualification. Taking its inner product with both sides of equation (5.1) gives

$$\sum_{j=1}^r \mu_j dh_j(\mathbf{0})\mathbf{v} = 0,$$

contradicting the assumption that $dh_j(\mathbf{0})\mathbf{v} < 0$ for all $1 \leq j \leq r$ and the fact that at least one $\mu_j > 0$. Thus, $\lambda_0 > 0$, and we can divide equation (5.1) by λ_0 . ■

Example 5.2.2 Application to an Inequality

Let us demonstrate the inequality

$$\frac{x_1^2 + x_2^2}{4} \leq e^{x_1 + x_2 - 2}$$

subject to the constraints $x_1 \geq 0$ and $x_2 \geq 0$ [69]. It suffices to show that the minimum of $f(\mathbf{x}) = -(x_1^2 + x_2^2)e^{-x_1 - x_2}$ is $-4e^{-2}$. According to Proposition 5.2.1 with $h_1(\mathbf{x}) = -x_1$ and $h_2(\mathbf{x}) = -x_2$, a minimum point entails the conditions

$$\begin{aligned} -\frac{\partial}{\partial x_1} f(\mathbf{x}) &= (2x_1 - x_1^2 - x_2^2)e^{-x_1 - x_2} = -\mu_1 \\ -\frac{\partial}{\partial x_2} f(\mathbf{x}) &= (2x_2 - x_1^2 - x_2^2)e^{-x_1 - x_2} = -\mu_2, \end{aligned}$$

where the multipliers μ_1 and μ_2 are nonnegative and satisfy $\mu_1 x_1 = 0$ and $\mu_2 x_2 = 0$. In this problem, the Mangasarian-Fromovitz constraint qualification is trivial to check using the vector $\mathbf{v} = \mathbf{1}$. If neither x_1 nor x_2 vanishes, then

$$2x_1 - x_1^2 - x_2^2 = 2x_2 - x_1^2 - x_2^2 = 0.$$

This forces $x_1 = x_2$ and $2x_1 - 2x_1^2 = 0$. It follows that $x_1 = x_2 = 1$, where $f(\mathbf{1}) = -2e^{-2}$. We can immediately eliminate the origin $\mathbf{0}$ from contention because $f(\mathbf{0}) = 0$. If $x_1 = 0$ and $x_2 > 0$, then $\mu_2 = 0$ and $2x_2 - x_2^2 = 0$. This implies that $x_2 = 2$ and $(0, 2)$ is a candidate minimum point. By symmetry, $(2, 0)$ is also a candidate minimum point. At these two boundary points, $f(2, 0) = f(0, 2) = -4e^{-2}$, and this verifies the claimed minimum value. ■

Example 5.2.3 *Application to Linear Programming*

The gradient of the Lagrangian

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathbf{z}^* \mathbf{x} + \sum_{i=1}^p \lambda_i \left(\sum_{j=1}^n v_{ij} x_j - d_i \right) - \sum_{j=1}^q \mu_j x_j$$

vanishes at the minimum of the linear function $f(\mathbf{x}) = \mathbf{z}^* \mathbf{x}$ subject to the constraints $\mathbf{V}\mathbf{x} = \mathbf{d}$ and $\mathbf{x} \geq \mathbf{0}$. Differentiating $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ with respect to x_j and setting the result equal to 0 gives $z_j + \sum_{i=1}^p \lambda_i v_{ij} - \mu_j = 0$. In vector notation this is just

$$\mathbf{z} + \mathbf{V}^* \boldsymbol{\lambda} - \boldsymbol{\mu} = \mathbf{0}$$

subject to the restrictions $\boldsymbol{\mu} \geq \mathbf{0}$ and $\boldsymbol{\mu}^* \mathbf{x} = 0$. We will revisit linear programming later and discuss an interior point algorithm and dual linear programs. ■

Example 5.2.4 *A Counterexample to the Multiplier Rule*

The Lagrange multiplier condition is necessary but not sufficient for a point to furnish a minimum. For example, consider the function $f(\mathbf{x}) = x_1^3 - x_2$ subject to the constraint $h(\mathbf{x}) = x_2 \leq 0$. The Lagrange multiplier condition

$$\nabla f(\mathbf{0}) = \begin{pmatrix} 0 \\ -1 \end{pmatrix} = -\nabla h(\mathbf{0})$$

holds, but the origin $\mathbf{0}$ fails to minimize $f(\mathbf{x})$. Indeed, the one-dimensional slice $x_1 \mapsto f(x_1, 0)$ has a saddle point at $x_1 = 0$. This function has no minimum subject to the inequality constraint. ■

Example 5.2.5 *Shadow Values*

In economic applications the objective function $f(\mathbf{x})$ is often viewed as a profit or a cost, and the constraints equation $g_i(\mathbf{x}) = 0$ is rephrased as $g_i(\mathbf{x}) = c_i$, where c_i is the amount of some available resource. In the absence of inequality constraints, the negative Lagrange multiplier $-\lambda_i$ is called a shadow value or price and measures the rate of change of the optimal value of $f(\mathbf{x})$ relative to c_i . Let $x(\mathbf{c})$ denote the optimal point as a function of the constraint vector $\mathbf{c} = (c_1, \dots, c_p)^*$. If we assume that $x(\mathbf{c})$ is differentiable, then we can multiply the equation

$$df[x(\mathbf{c})] + \sum_{i=1}^p \lambda_i dg_i[x(\mathbf{c})] = \mathbf{0}^*$$

on the right by $dx(\mathbf{c})$ and recover via the chain rule an equation relating the differentials of $f[x(\mathbf{c})]$ and the $g_i[x(\mathbf{c})]$ with respect to \mathbf{c} . Because it is obvious that

$$\frac{\partial}{\partial c_j} g_i[x(\mathbf{c})] = 1_{\{j=i\}},$$

it follows that

$$\frac{\partial}{\partial c_j} f[x(\mathbf{c})] + \lambda_j = 0.$$

Of course, this result is valid generally and transcends its narrow economic origin. ■

Example 5.2.6 *Quadratic Programming with Equality Constraints*

To minimize the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ subject to the linear equality constraints $\mathbf{V} \mathbf{x} = \mathbf{d}$, we introduce the Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + \sum_{i=1}^p \lambda_i [\mathbf{v}_i^* \mathbf{x} - d_i] \\ &= \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + \boldsymbol{\lambda}^* (\mathbf{V} \mathbf{x} - \mathbf{d}). \end{aligned}$$

A stationary point of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is determined by the equations

$$\begin{aligned} \mathbf{A} \mathbf{x} + \mathbf{b} + \mathbf{V}^* \boldsymbol{\lambda} &= \mathbf{0} \\ \mathbf{V} \mathbf{x} &= \mathbf{d}, \end{aligned}$$

whose formal solution amounts to

$$\begin{pmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \end{pmatrix} = \begin{pmatrix} \mathbf{A} & \mathbf{V}^* \\ \mathbf{V} & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} -\mathbf{b} \\ \mathbf{d} \end{pmatrix}.$$

The next proposition shows that the indicated matrix inverse exists when \mathbf{A} is positive definite. ■

Proposition 5.2.2 *Let \mathbf{A} be an $n \times n$ positive definite matrix and \mathbf{V} be a $p \times n$ matrix. Then the matrix*

$$\mathbf{M} = \begin{pmatrix} \mathbf{A} & \mathbf{V}^* \\ \mathbf{V} & \mathbf{0} \end{pmatrix}$$

is invertible if and only if \mathbf{V} has linearly independent rows $\mathbf{v}_1^, \dots, \mathbf{v}_p^*$. When this condition holds, \mathbf{M} has inverse*

$$\mathbf{M}^{-1} = \begin{pmatrix} \mathbf{A}^{-1} - \mathbf{A}^{-1} \mathbf{V}^* (\mathbf{V} \mathbf{A}^{-1} \mathbf{V}^*)^{-1} \mathbf{V} \mathbf{A}^{-1} & \mathbf{A}^{-1} \mathbf{V}^* (\mathbf{V} \mathbf{A}^{-1} \mathbf{V}^*)^{-1} \\ (\mathbf{V} \mathbf{A}^{-1} \mathbf{V}^*)^{-1} \mathbf{V} \mathbf{A}^{-1} & -(\mathbf{V} \mathbf{A}^{-1} \mathbf{V}^*)^{-1} \end{pmatrix}.$$

Proof: We first show that the symmetric matrix \mathbf{M} is invertible with the specified inverse if and only if $(\mathbf{V} \mathbf{A}^{-1} \mathbf{V}^*)^{-1}$ exists. If \mathbf{M}^{-1} exists, it is necessarily symmetric. Indeed, taking the transpose of $\mathbf{M} \mathbf{M}^{-1} = \mathbf{I}$ gives $(\mathbf{M}^{-1})^* \mathbf{M} = \mathbf{I}$. Suppose \mathbf{M}^{-1} has block form $\begin{pmatrix} \mathbf{B} & \mathbf{C}^* \\ \mathbf{C} & \mathbf{D} \end{pmatrix}$. Then the identity

$$\begin{pmatrix} \mathbf{A} & \mathbf{V}^* \\ \mathbf{V} & \mathbf{0} \end{pmatrix} \begin{pmatrix} \mathbf{B} & \mathbf{C}^* \\ \mathbf{C} & \mathbf{D} \end{pmatrix} = \begin{pmatrix} \mathbf{I}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_p \end{pmatrix}$$

implies that $\mathbf{V}\mathbf{C}^* = \mathbf{I}_p$ and $\mathbf{A}\mathbf{C}^* + \mathbf{V}^*\mathbf{D} = \mathbf{0}$. Multiplying the last equality by $\mathbf{V}\mathbf{A}^{-1}$ gives $\mathbf{I}_p = -\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*\mathbf{D}$. Thus, $(\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*)^{-1}$ exists. Conversely, if $(\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*)^{-1}$ exists, then one can check by direct multiplication that \mathbf{M} has the claimed inverse.

If $(\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*)^{-1}$ exists, then \mathbf{V} must have full row rank p . Conversely, if \mathbf{V} has full row rank p , take any nontrivial $\mathbf{u} \in \mathbb{R}^p$. The fact

$$\mathbf{u}^*\mathbf{V} = u_1\mathbf{v}_1^* + \cdots + u_p\mathbf{v}_p^* \neq \mathbf{0}^*$$

and the positive definiteness of \mathbf{A} imply $\mathbf{u}^*\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*\mathbf{u} > 0$. Thus, $\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*$ is positive definite and invertible. ■

Example 5.2.7 *Smallest Matrix Subject to Secant Conditions*

In some situations covered by Example 5.2.6, the answer can be radically simplified. Consider the problem of minimizing the Frobenius norm of a matrix \mathbf{M} subject to the linear constraints $\mathbf{M}\mathbf{u}_i = \mathbf{v}_i$ for $i = 1, \dots, q$. It is helpful to rewrite the constraints in matrix form as $\mathbf{M}\mathbf{U} = \mathbf{V}$ for $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_q)$ and $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_q)$. Provided \mathbf{U} has full column rank q , the minimum of the squared norm $\|\mathbf{M}\|_F^2$ subject to the constraints is attained by the choice $\mathbf{M} = \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*$. We can prove this assertion by taking the partial derivative of the Lagrangian

$$\begin{aligned} \mathcal{L} &= \frac{1}{2}\|\mathbf{M}\|_F^2 + \sum_i \sum_k \lambda_{ik} \left(\sum_j m_{ij}u_{jk} - v_{ik} \right) \\ &= \frac{1}{2} \sum_i \sum_j m_{ij}^2 + \sum_i \sum_k \lambda_{ik} \left(\sum_j m_{ij}u_{jk} - v_{ik} \right) \end{aligned}$$

with respect to m_{ij} and equating it to 0. This gives the Lagrange multiplier equation

$$0 = m_{ij} + \sum_k \lambda_{ik}u_{jk},$$

which we collectively express in matrix notation as $\mathbf{0} = \mathbf{M} + \mathbf{\Lambda}\mathbf{U}^*$. This equation and the constraint equation $\mathbf{M}\mathbf{U} = \mathbf{V}$ uniquely determine the minimum of the objective function. Indeed, straightforward substitution shows that $\mathbf{M} = \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*$ and $\mathbf{\Lambda} = -\mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}$ constitute the solution. This result will come in handy later when we discuss accelerating various algorithms. ■

5.3 Constraint Qualification

Fortunately, the Mangasarian-Fromovitz constraint qualification is a consequence of a stronger condition suggested by Kuhn and Tucker. This alternative condition requires the differentials $dg_i(\mathbf{x})$ of the equality constraints and the differentials $dh_j(\mathbf{x})$ of the active inequality constraints to

be collectively linearly independent at \mathbf{x} . When the Kuhn-Tucker condition is true, we can find an $n \times n$ invertible matrix \mathbf{A} whose first p rows are the vectors $dg_i(\mathbf{x})$ and whose next r rows are the vectors $dh_1(\mathbf{x}), \dots, dh_r(\mathbf{x})$ corresponding to the active inequality constraints. For example, we can choose the last $n - p - r$ rows of \mathbf{A} to be a basis for the orthogonal complement of the subspace spanned by the first $p + r$ rows of \mathbf{A} . Given that \mathbf{A} is invertible, there certainly exists a vector $\mathbf{v} \in \mathbb{R}^n$ with

$$\mathbf{A}\mathbf{v} = - \begin{pmatrix} \mathbf{0} \\ \mathbf{1} \\ \mathbf{0} \end{pmatrix},$$

where the column vectors on the right side of this equality have p , r , and $n - p - r$ rows, respectively. Clearly, the vector \mathbf{v} satisfies the Mangasarian-Fromovitz constraint qualification.

Even if the Kuhn-Tucker condition fails, there is still a chance that the constraint qualification holds. Let \mathbf{G} be the $p \times n$ matrix with rows $dg_i(\mathbf{x})$. Assume that \mathbf{G} has full rank. Taking the orthogonal complement of the subspace spanned by the p rows of \mathbf{G} , we can construct an $n \times (n - p)$ matrix \mathbf{K} of full rank whose columns are orthogonal to the rows of \mathbf{G} . This fact can be expressed as $\mathbf{G}\mathbf{K} = \mathbf{0}$ and implies that the image of the linear transformation \mathbf{K} from \mathbb{R}^{n-p} to \mathbb{R}^n is the kernel or null space of \mathbf{G} . In other words, any \mathbf{v} satisfying $\mathbf{G}\mathbf{v} = \mathbf{0}$ can be expressed as $\mathbf{v} = \mathbf{K}\mathbf{u}$ and vice versa. If we let $\mathbf{z}_j^* = dh_j(\mathbf{x})\mathbf{K}$ for each of the active inequality constraints, then the Mangasarian-Fromovitz constraint qualification is equivalent to the existence of a vector $\mathbf{u} \in \mathbb{R}^{n-p}$ satisfying $\mathbf{z}_j^*\mathbf{u} < 0$ for all $1 \leq j \leq r$. If the number of equality constraints $p = 0$, then we take $\mathbf{z}_j^* = dh_j(\mathbf{x})$.

The next proposition paves the way for proving a necessary and sufficient condition for this equality-free form of the Mangasarian-Fromovitz constraint qualification. The result described by the proposition is of independent interest; its proof illustrates the fact that adding small penalties as opposed to large penalties is sometimes helpful [84].

Proposition 5.3.1 (Ekeland) *Suppose the real-valued function $f(\mathbf{x})$ is defined and differentiable on \mathbb{R}^n . If $f(\mathbf{x})$ is bounded below, then there are points where $\|\nabla f(\mathbf{x})\|$ is arbitrarily close to 0.*

Proof: Take any small $\epsilon > 0$ and define the continuous function

$$f_\epsilon(\mathbf{x}) = f(\mathbf{x}) + \epsilon\|\mathbf{x}\|.$$

In view of the boundedness condition, for any \mathbf{y} the set

$$\{\mathbf{x} : f_\epsilon(\mathbf{x}) \leq f_\epsilon(\mathbf{y})\}$$

is compact. Hence, Proposition 2.5.4 implies that $f_\epsilon(\mathbf{x})$ has a global minimum \mathbf{z} depending on ϵ . We now prove that $\mathbf{v} = \nabla f(\mathbf{z})$ satisfies $\|\mathbf{v}\| \leq \epsilon$. If the opposite is true, then the limit relation

$$\begin{aligned} \lim_{t \downarrow 0} \frac{f(\mathbf{z} - t\mathbf{v}) - f(\mathbf{z})}{t} &= -df(\mathbf{z})\mathbf{v} \\ &= -\|\mathbf{v}\|^2 \\ &< -\epsilon\|\mathbf{v}\|, \end{aligned}$$

the choice of \mathbf{z} , and the triangle inequality together entail

$$\begin{aligned} -t\epsilon\|\mathbf{v}\| &> f(\mathbf{z} - t\mathbf{v}) - f(\mathbf{z}) \\ &= f_\epsilon(\mathbf{z} - t\mathbf{v}) - f_\epsilon(\mathbf{z}) + \epsilon\|\mathbf{z}\| - \epsilon\|\mathbf{z} - t\mathbf{v}\| \\ &\geq \epsilon\|\mathbf{z}\| - \epsilon\|\mathbf{z} - t\mathbf{v}\| \\ &\geq -\epsilon t\|\mathbf{v}\| \end{aligned}$$

for sufficiently small $t > 0$. This contradiction implies that $\|\mathbf{v}\| \leq \epsilon$. ■

We are now in position to characterize the Mangasarian-Fromovitz constraint qualification in the absence of equality constraints. As the next proposition indicates, the constraint qualification holds if and only if the convex set generated by the active inequality constraints does not contain the origin. The geometric nature of this result will be clearer after we consider convex sets in the next chapter.

Proposition 5.3.2 (Gordon) *Given r vectors $\mathbf{z}_1, \dots, \mathbf{z}_r$ in \mathbb{R}^n , define the function $f(\mathbf{x}) = \ln \left[\sum_{j=1}^r \exp(\mathbf{z}_j^* \mathbf{x}) \right]$. Then the following three conditions are logically equivalent:*

- (a) *The function $f(\mathbf{x})$ is bounded below on \mathbb{R}^n .*
- (b) *There are nonnegative constants μ_1, \dots, μ_r such that*

$$\sum_{i=1}^r \mu_i \mathbf{z}_i = \mathbf{0}, \quad \sum_{i=1}^r \mu_i = 1.$$

- (c) *There is no vector \mathbf{u} such that $\mathbf{z}_j^* \mathbf{u} < 0$ for all j .*

Proof: It is trivial to check that (b) implies (c) and (c) implies (a). To demonstrate that (a) implies (b), first observe that

$$\nabla f(\mathbf{x}) = \frac{1}{\sum_{i=1}^r e^{\mathbf{z}_i^* \mathbf{x}}} \sum_{j=1}^r e^{\mathbf{z}_j^* \mathbf{x}} \mathbf{z}_j.$$

According to Proposition 5.3.1, it is possible to choose for each k a point \mathbf{u}_k at which

$$\|\nabla f(\mathbf{u}_k)\| = \left\| \sum_{j=1}^r \mu_{kj} \mathbf{z}_j \right\| \leq \frac{1}{k},$$

where

$$\mu_{kj} = \frac{e^{\mathbf{z}_j^* \mathbf{u}_k}}{\sum_{i=1}^r e^{\mathbf{z}_i^* \mathbf{u}_k}}.$$

Because the μ_{kj} form a vector $\boldsymbol{\mu}_k$ with nonnegative components summing to 1, it is possible to find a subsequence of the sequence $\boldsymbol{\mu}_k$ that converges to a vector $\boldsymbol{\mu}$ having the same properties. This vector satisfies the requirements of condition (b). ■

When the equality constraints in nonlinear programming are affine, the forgoing discussion is easy to summarize. One first eliminates the equality constraints via the reparameterization $\mathbf{x} = \mathbf{K}\mathbf{y}$. Then one defines the differentials $\mathbf{z}_j^* = dh_j(\mathbf{x})\mathbf{K}$ of the inequality constraints $h_j(\mathbf{K}\mathbf{y})$ in the new parameterization. The Mangasarian-Fromovitz constraint qualification postulates the existence of a vector \mathbf{u} with $\mathbf{u}^* \mathbf{z}_j < 0$ for all active constraints $h_j(\mathbf{x})$. The existence of \mathbf{u} rules out the possibility that $\mathbf{0}$ is a convex combination of the vectors \mathbf{z}_j . Hence, the Lagrange multiplier λ_0 of $f(\mathbf{x})$ in the \mathbf{y} parameterization cannot equal 0.

5.4 Taylor-Made Higher-Order Differentials

Roughly speaking, the higher-order differentials of a function $f(\mathbf{x})$ on \mathbb{R}^n arise from its multiple partial derivatives. Thus, $\partial_{ij}^2 f(\mathbf{x})$ contributes to a second differential and $\partial_{ijk}^3 f(\mathbf{x})$ to a third differential. Equality of mixed partial derivatives involves simple identities such as $\partial_{ij}^2 f(\mathbf{x}) = \partial_{ji}^2 f(\mathbf{x})$ and $\partial_{ijk}^3 f(\mathbf{x}) = \partial_{kji}^3 f(\mathbf{x})$. The right notation efficiently exposes these symmetries. Let \mathbf{j} be a vector whose m components j_i are drawn from the set $\{1, \dots, n\}$. In this notation we write a mixed partial derivative of the function $f(\mathbf{x})$ as

$$\partial_{\mathbf{j}}^m f(\mathbf{x}) = \partial_{j_1} \cdots \partial_{j_m} f(\mathbf{x}).$$

If the partial derivatives of $f(\mathbf{x})$ of order m are continuous at $\mathbf{x} \in \mathbb{R}^n$, then Proposition 4.3.1 shows that equality of mixed partials holds and the order of the components of \mathbf{j} is irrelevant. The multinomial coefficient

$$\binom{m}{\mathbf{k}} = \binom{m}{k_1 \dots k_n}$$

counts the number of mixed partial derivatives of order m in which the partial differential operator ∂_i appears k_i times.

Suppose $f(\mathbf{y})$ is a real-valued function possessing all partial derivatives of order p or less in a neighborhood of the point \mathbf{x} . The first-order Taylor expansion

$$f(\mathbf{y}) = f(\mathbf{x}) + \sum_{i=1}^n \int_0^1 \partial_i f[t\mathbf{y} + (1-t)\mathbf{x}] dt (y_i - x_i)$$

follows directly from the fundamental theorem of calculus and the chain rule. To pass to the second-order Taylor expansion, one integrates by parts and replaces the integral

$$\int_0^1 \partial_i f[t\mathbf{y} + (1-t)\mathbf{x}] dt (y_i - x_i)$$

by

$$\partial_i f(\mathbf{x})(y_i - x_i) + \sum_{j=1}^n \int_0^1 \partial_j \partial_i f[t\mathbf{y} + (1-t)\mathbf{x}](1-t) dt (y_i - x_i)(y_j - x_j).$$

Repeated integration by parts, differentiating $f(\mathbf{x})$ and integrating successive powers of $(1-t)$, ultimately leads to the Taylor expansion

$$\begin{aligned} f(\mathbf{y}) &= \sum_{m=0}^{p-1} \frac{1}{m!} d^m f(\mathbf{x})[(\mathbf{y} - \mathbf{x})^m] + R(\mathbf{y}, \mathbf{x}) \\ R(\mathbf{y}, \mathbf{x}) &= \frac{1}{(p-1)!} \int_0^1 d^p f[t\mathbf{y} + (1-t)\mathbf{x}][(\mathbf{y} - \mathbf{x})^p](1-t)^{p-1} dt \end{aligned} \quad (5.7)$$

of order p . Here we employ the notation of multilinear maps sketched in Example 2.5.10. The abstract entity $d^m f(\mathbf{x})[(\mathbf{y} - \mathbf{x})^m]$ is a m -linear form evaluated along its diagonal whose coefficients relative to the standard basis of \mathbb{R}^n reduce to the mixed partial derivatives $\partial_j^m f(\mathbf{x})$ of order m . Equality of mixed partials makes $d^m f(\mathbf{x})$ a symmetric m -linear form. All of this sounds complicated, but it simply amounts to local approximation of $f(\mathbf{y})$ by a polynomial in the components of the difference vector $\mathbf{y} - \mathbf{x}$. Our derivation of the Taylor expansion (5.7) remains valid for $f(\mathbf{y})$ vector or matrix valued if we operate entry by entry.

The remainder $R(\mathbf{y}, \mathbf{x})$ in the expansion (5.7) can be recast in two useful ways. Assuming the partial derivatives of $f(\mathbf{x})$ of order p and less are continuous at \mathbf{x} , it is obvious that

$$s^p(\mathbf{y}, \mathbf{x})[\mathbf{u}_1, \dots, \mathbf{u}_p] = p \int_0^1 d^p f[t\mathbf{y} + (1-t)\mathbf{x}][\mathbf{u}_1, \dots, \mathbf{u}_p](1-t)^{p-1} dt$$

is an p -linear map with $\lim_{\mathbf{y} \rightarrow \mathbf{x}} s^p(\mathbf{y}, \mathbf{x}) = s^p(\mathbf{x}, \mathbf{x}) = d^p f(\mathbf{x})$. This suggests that the expansion

$$f(\mathbf{y}) = \sum_{m=0}^{p-1} \frac{1}{m!} d^m f(\mathbf{x})[(\mathbf{y} - \mathbf{x})^m] + \frac{1}{p!} s^p(\mathbf{y}, \mathbf{x})[(\mathbf{y} - \mathbf{x})^p] \quad (5.8)$$

be taken as the definition of Carathéodory differentiability of order p , with the understanding that the slope $s^p(\mathbf{y}, \mathbf{x})$ tends to $s^p(\mathbf{x}, \mathbf{x}) = d^p f(\mathbf{x})$ as \mathbf{y} tends to \mathbf{x} . There is no harm in assuming that $s^p(\mathbf{y}, \mathbf{x})$ and $d^p f(\mathbf{x})$ are

symmetric multilinear maps since they operate only on diagonal arguments. Indeed, any multilinear map $M[\mathbf{u}_1, \dots, \mathbf{u}_p]$ agrees with its symmetrization

$$M_{\text{Sym}}[\mathbf{u}_1, \dots, \mathbf{u}_p] = \frac{1}{p!} \sum_{\sigma} M[\mathbf{u}_{\sigma_1}, \dots, \mathbf{u}_{\sigma_p}]$$

along its diagonal $\mathbf{u}_1 = \mathbf{u}_2 = \dots = \mathbf{u}_p$. Here the sum ranges over all permutations σ of $\{1, \dots, p\}$.

Alternatively, if we let $r^p(\mathbf{y}, \mathbf{x}, t) = d^p f[t\mathbf{y} + (1-t)\mathbf{x}] - d^p f(\mathbf{x})$, then the identity

$$R(\mathbf{y}, \mathbf{x}) = \frac{1}{(p-1)!} \int_0^1 r^p(\mathbf{y}, \mathbf{x}, t)[(\mathbf{y} - \mathbf{x})^p](1-t)^{p-1} dt + \frac{1}{p!} d^p f(\mathbf{x})[(\mathbf{y} - \mathbf{x})^p]$$

and the inequality

$$\|r^p(\mathbf{y}, \mathbf{x}, t)[(\mathbf{y} - \mathbf{x})^p]\| \leq \|r^p(\mathbf{y}, \mathbf{x}, t)\| \cdot \|\mathbf{y} - \mathbf{x}\|^p$$

suggest that the expansion

$$f(\mathbf{y}) = \sum_{m=0}^p \frac{1}{m!} d^m f(\mathbf{x})[(\mathbf{y} - \mathbf{x})^m] + o(\|\mathbf{y} - \mathbf{x}\|^p)$$

be taken as the definition of Fréchet differentiability of order p . These two axiomatic definitions of higher-order differentials are logically equivalent. Given Carathéodory's definition, we set $d^p f(\mathbf{x}) = s^p(\mathbf{x}, \mathbf{x})$ and note that

$$\| [s^p(\mathbf{y}, \mathbf{x}) - s^p(\mathbf{x}, \mathbf{x})][(\mathbf{y} - \mathbf{x})^p] \| \leq \|s^p(\mathbf{y}, \mathbf{x}) - s^p(\mathbf{x}, \mathbf{x})\| \|\mathbf{y} - \mathbf{x}\|^p.$$

Fréchet's definition follows directly. Proof of the converse is more subtle, and we omit it. Both definitions require the existence of all partial derivatives of $f(\mathbf{y})$ of order $p - 1$ and less in a neighborhood of \mathbf{x} .

Perhaps just as important as the equivalence of the two definition of differentiability is the following simple result.

Proposition 5.4.1 *The symmetric multilinear maps $d^m f(\mathbf{x})$ and $s^p(\mathbf{x}, \mathbf{x})$ appearing in the Taylor expansion (5.8) are unique. The slope function $s^p(\mathbf{y}, \mathbf{x})$ is not unique as noted in Chap. 4.*

Proof: Abbreviate $d^m f(\mathbf{x}) = M^m$. The constant M^0 equals $\lim_{\mathbf{y} \rightarrow \mathbf{x}} f(\mathbf{y})$. Consider a second expansion of $f(\mathbf{y})$ with N^m replacing M^m and $t^p(\mathbf{y}, \mathbf{x})$ replacing $s^p(\mathbf{y}, \mathbf{x})$. Suppose by induction that $M^k = N^k$ for $k < m < p$, and let \mathbf{y}_j be the sequence $\mathbf{x} + \frac{1}{j}\mathbf{u}$ for some fixed vector \mathbf{u} . Taking limits on j in the equality

$$\begin{aligned} \mathbf{0} &= \frac{1}{m!} (M^m - N^m)[\mathbf{u}^m] + \sum_{k=m+1}^{p-1} \frac{j^{m-k}}{k!} (M^k - N^k)[\mathbf{u}^k] \\ &\quad + j^{m-p} [s^p(\mathbf{y}_j, \mathbf{x}) - t^p(\mathbf{y}_j, \mathbf{x})][\mathbf{u}^p] \end{aligned}$$

gives $\mathbf{M}^m[\mathbf{u}^m] = \mathbf{N}^m[\mathbf{u}^m]$. If we now interpret $\mathbf{M}^m[\mathbf{u}^m]$ and $\mathbf{N}^m[\mathbf{u}^m]$ as polynomials in the entries of \mathbf{u} , then their coefficients must coincide. In view of the symmetry of \mathbf{M}^m and \mathbf{N}^m , the coefficients $\mathbf{M}^m[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_m}]$ and $\mathbf{N}^m[\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_m}]$ therefore also coincide for all choices $\mathbf{e}_{i_1}, \dots, \mathbf{e}_{i_m}$. This argument advances the induction one step. The proof of the equality $s^p(\mathbf{x}, \mathbf{x}) = t^p(\mathbf{x}, \mathbf{x})$ follows essentially the same pattern. ■

Example 5.4.1 *Second Differential of a Bilinear Map*

Consider a bilinear map $f(\mathbf{x}) = M(\mathbf{x}_1, \mathbf{x}_2)$ of the concatenated vector \mathbf{x} with left block \mathbf{x}_1 and right block \mathbf{x}_2 . If we take $\mathbf{u}_i = \mathbf{y}_i - \mathbf{x}_i$, then the expansion

$$\begin{aligned} M(\mathbf{y}_1, \mathbf{y}_2) &= M(\mathbf{x}_1 + \mathbf{u}_1, \mathbf{x}_2 + \mathbf{u}_2) \\ &= M(\mathbf{x}_1, \mathbf{x}_2) + M(\mathbf{x}_1, \mathbf{u}_2) + M(\mathbf{u}_1, \mathbf{x}_2) + M(\mathbf{u}_1, \mathbf{u}_2) \end{aligned}$$

identifies the bilinear map $2M(\mathbf{u}_1, \mathbf{u}_2)$ as the second differential of $f(\mathbf{x})$. Similar expansions hold for higher-order multilinear maps. ■

The most important applications in optimization theory of higher-order differentials involve scalar-valued functions $f(\mathbf{x})$ and their first and second differentials $df(\mathbf{x})$ and $d^2f(\mathbf{x})$. We retain our convention that the gradient $\nabla f(\mathbf{x})$ equals the transpose of the differential $df(\mathbf{x})$. The coefficients of the linear form $df(\mathbf{x})[\mathbf{u}]$ are the first partials $\partial_i f(\mathbf{x})$. The second-order Taylor expansion of $f(\mathbf{y})$ around \mathbf{x} reads

$$f(\mathbf{y}) = f(\mathbf{x}) + df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^* s^2(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}). \quad (5.9)$$

Here the second-order slope $s^2(\mathbf{y}, \mathbf{x})$ has limit $d^2f(\mathbf{x})$; this Hessian matrix incorporates the second partials $\partial_{ij}^2 f(\mathbf{x})$. Consequently, one is justified in viewing $d^2f(\mathbf{x})$ as the differential of $\nabla f(\mathbf{x})$.

Example 5.4.2 *Second Differential of a Quadratic Function*

Consider the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ defined by an $n \times n$ symmetric matrix \mathbf{A} and an $n \times 1$ vector \mathbf{b} . The gradient of $f(\mathbf{y})$ is $\nabla f(\mathbf{y}) = \mathbf{A} \mathbf{y} + \mathbf{b}$. The second differential emerges from the slope equation

$$\nabla f(\mathbf{y}) - \nabla f(\mathbf{x}) = \mathbf{A}(\mathbf{y} - \mathbf{x})$$

or the simple calculation $d^2f(\mathbf{x}) = d\nabla f(\mathbf{x}) = (d\mathbf{A})\mathbf{y} + \mathbf{A}d\mathbf{y} + d\mathbf{b} = \mathbf{A}$. Either perspective is consistent with the exact expansion

$$f(\mathbf{y}) = f(\mathbf{x}) + (\mathbf{A} \mathbf{x} + \mathbf{b})^*(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^* \mathbf{A}(\mathbf{y} - \mathbf{x}).$$

Uniqueness of the differentials $df(\mathbf{x})$ and $d^2f(\mathbf{x})$ is guaranteed by Proposition 5.4.1. ■

Example 5.4.3 *A Pathological Example*

On \mathbb{R}^2 define the indicator function $1_Q(\mathbf{x})$ to be 1 when both coordinates x_1 and x_2 are rational numbers and to be 0 otherwise. The function

$$f(\mathbf{x}) = 1_Q(\mathbf{x})(x_1^3 + x_2^3)$$

is differentiable at the origin $\mathbf{0}$ but nowhere else on \mathbb{R}^2 . The differential $df(\mathbf{0}) = \mathbf{0}$ together with the choice

$$s^2(\mathbf{y}, \mathbf{0}) = 1_Q(\mathbf{y}) \begin{pmatrix} 2y_1 & 0 \\ 0 & 2y_2 \end{pmatrix}.$$

satisfy expansion (5.9). Furthermore, $s^2(\mathbf{y}, \mathbf{0})$ tends to the $\mathbf{0}$ matrix as \mathbf{y} tends to the origin. However, claiming that $f(\mathbf{x})$ is twice differentiable at the origin seems extreme. To eliminate such pathologies, we demand that $f(\mathbf{x})$ be differentiable in a neighborhood of a point before we allow it to be twice differentiable at the point. ■

Example 5.4.4 *Second Differential of the Inverse Polar Transformation*

Continuing Example 4.6.1, it is straightforward to calculate that the inverse polar transformation

$$g \left[\begin{pmatrix} r \\ \theta \end{pmatrix} \right] = \begin{pmatrix} r \cos \theta \\ r \sin \theta \end{pmatrix}$$

has second differential

$$d^2g \left[\begin{pmatrix} r \\ \theta \end{pmatrix} \right] = \begin{pmatrix} d^2 r \cos \theta \\ d^2 r \sin \theta \end{pmatrix} = \begin{pmatrix} 0 & -\sin \theta \\ -\sin \theta & -r \cos \theta \\ 0 & \cos \theta \\ \cos \theta & -r \sin \theta \end{pmatrix}.$$

Here we stack the differentials corresponding to the two components for easy visualization. ■

Example 5.4.5 *First and Second Differentials of the Inverse of a Matrix*

Taylor expansions take different shapes depending on how one displays the results. For instance, the identity

$$\mathbf{Y}^{-1} - \mathbf{X}^{-1} = -\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{Y}^{-1}$$

is a disguised slope expansion of the inverse of a matrix as a function of its entries. The corresponding first differential is the linear map

$$\mathbf{U} \mapsto -\mathbf{X}^{-1}\mathbf{U}\mathbf{X}^{-1}.$$

The more elaborate second-order expansion

$$\begin{aligned} \mathbf{Y}^{-1} - \mathbf{X}^{-1} &= -\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{X}^{-1} + \frac{1}{2}\mathbf{Y}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{X}^{-1} \\ &\quad + \frac{1}{2}\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{X}^{-1}(\mathbf{Y} - \mathbf{X})\mathbf{Y}^{-1} \end{aligned}$$

shows that the second differential is the bilinear map

$$(\mathbf{U}, \mathbf{V}) \mapsto \mathbf{X}^{-1}\mathbf{U}\mathbf{X}^{-1}\mathbf{V}\mathbf{X}^{-1} + \mathbf{X}^{-1}\mathbf{V}\mathbf{X}^{-1}\mathbf{U}\mathbf{X}^{-1}$$

evaluated along $\mathbf{U} = \mathbf{V}$. Despite the odd appearance of the second differential, it clearly fulfills its quadratic approximation responsibility. ■

The rules for calculating second differentials are naturally more complicated than those for calculating first differentials.

Proposition 5.4.2 *If the two functions $f(\mathbf{x})$ and $g(\mathbf{x})$ map the open set $S \subset \mathbb{R}^p$ twice differentiably into \mathbb{R}^q , then the following functional combinations are twice differentiable and have the indicated second differentials:*

(a) For all constants α and β ,

$$d^2[\alpha f(\mathbf{x}) + \beta g(\mathbf{x})] = \alpha d^2 f(\mathbf{x}) + \beta d^2 g(\mathbf{x}).$$

(b) The inner product $f(\mathbf{x})^* g(\mathbf{x})$ satisfies

$$\begin{aligned} d^2[f(\mathbf{x})^* g(\mathbf{x})] &= \sum_{i=1}^q [f_i(\mathbf{x})d^2 g_i(\mathbf{x}) + g_i(\mathbf{x})d^2 f_i(\mathbf{x})] \\ &\quad + \sum_{i=1}^q [\nabla f_i(\mathbf{x})dg_i(\mathbf{x}) + \nabla g_i(\mathbf{x})df_i(\mathbf{x})]. \end{aligned}$$

(c) For $q = 1$ and $f(\mathbf{x}) \neq 0$,

$$d^2[f(\mathbf{x})^{-1}] = 2f(\mathbf{x})^{-3}\nabla f(\mathbf{x})df(\mathbf{x}) - f(\mathbf{x})^{-2}d^2 f(\mathbf{x}).$$

Proof: Rule (a) follows directly from the linearity implicit in formula (5.9) covering the scalar case. For rule (b) it also suffices to consider the scalar case in view of rule (a). Applying the sum and product rules of differentiation to the gradient of $f(\mathbf{x})g(\mathbf{x})$ gives

$$\begin{aligned} d[\nabla g(\mathbf{x})f(\mathbf{x}) + \nabla f(\mathbf{x})g(\mathbf{x})] &= d^2 g(\mathbf{x})f(\mathbf{x}) + \nabla g(\mathbf{x})df(\mathbf{x}) \\ &\quad + d^2 f(\mathbf{x})g(\mathbf{x}) + \nabla f(\mathbf{x})dg(\mathbf{x}). \end{aligned}$$

To verify rule (c), we apply the product, quotient, and chain rules to the gradient of $f(\mathbf{y})^{-1}$ identified in the proof of Proposition 4.4.2. This gives

$$-d\left[\nabla f(\mathbf{x})\frac{1}{f(\mathbf{x})^2}\right] = -d^2 f(\mathbf{x})\frac{1}{f(\mathbf{x})^2} + \nabla f(\mathbf{x})\frac{2}{f(\mathbf{x})^3}df(\mathbf{x}).$$

Note the care exercised here in ordering the different factors prior to differentiation. ■

The chain rule for second differentials also is more complex.

Proposition 5.4.3 *Suppose $f(\mathbf{x})$ maps the open set $S \subset \mathbb{R}^p$ twice differentiable into \mathbb{R}^q and $g(\mathbf{y})$ maps the open set $T \subset \mathbb{R}^q$ twice differentiable into \mathbb{R}^r . If the image $f(S)$ is contained in T , then the composite function $h(\mathbf{x}) = g \circ f(\mathbf{x})$ is twice differentiable with second partial derivatives*

$$\partial_{kl}^2 h_m(\mathbf{x}) = \sum_{i=1}^q \partial_i g_m \circ f(\mathbf{x}) \partial_{kl}^2 f_i(\mathbf{x}) + \sum_{i=1}^q \sum_{j=1}^q \partial_k f_i(\mathbf{x}) \partial_{ij}^2 g_m \circ f(\mathbf{x}) \partial_l f_j(\mathbf{x}).$$

Proof: It suffices to prove the result when $r = 1$ and $g(\mathbf{x})$ is scalar valued. The function $h(\mathbf{x})$ has first differential $dh(\mathbf{x}) = (dg) \circ f(\mathbf{x})df(\mathbf{x})$ and gradient $df(\mathbf{x})^*(\nabla g) \circ f(\mathbf{x})$. The matrix transpose, chain, and product rules of differentiation derived in Examples 4.7.1, 4.7.2, and 4.7.3 show that $\nabla h(\mathbf{x})$ has differential components

$$\begin{aligned} \partial_k [df(\mathbf{x})^* \nabla g \circ f(\mathbf{x})] &= [\partial_k df(\mathbf{x})]^* \nabla g \circ f(\mathbf{x}) \\ &\quad + df(\mathbf{x})^* \sum_{j=1}^q (\partial_j \nabla g) \circ f(\mathbf{x}) \partial_k f_j(\mathbf{x}). \end{aligned}$$

Alternatively, one can calculate the conventional way with explicit partial derivatives and easily verify the claimed formula for $d^2 h(\mathbf{x})$. ■

5.5 Applications of Second Differentials

Our first proposition, mentioned informally in Chap. 1, emphasizes the importance of second differentials in optimization theory.

Proposition 5.5.1 *Consider a real-valued function $f(\mathbf{y})$ with domain an open set $U \subset \mathbb{R}^p$. If $f(\mathbf{y})$ has a local minimum at \mathbf{x} and is twice differentiable there, then the second differential $d^2 f(\mathbf{x})$ is positive semidefinite. Conversely, if \mathbf{x} is a stationary point of $f(\mathbf{y})$ and $d^2 f(\mathbf{x})$ is positive definite, then \mathbf{x} is a local minimum of $f(\mathbf{y})$. Similar statements hold for local maxima if we replace the modifiers positive semidefinite and positive definite by the modifiers negative semidefinite and negative definite. Finally, if \mathbf{x} is a stationary point of $f(\mathbf{y})$ and $d^2 f(\mathbf{x})$ possesses both positive and negative eigenvalues, then \mathbf{x} is neither a local minimum nor a local maximum of $f(\mathbf{y})$.*

Proof: Suppose \mathbf{x} provides a local minimum of $f(\mathbf{y})$. For any unit vector \mathbf{v} and $t > 0$ sufficiently small, the point $\mathbf{y} = \mathbf{x} + t\mathbf{v}$ belongs to U and satisfies $f(\mathbf{y}) \geq f(\mathbf{x})$. If we divide the expansion

$$\begin{aligned} 0 &\leq f(\mathbf{y}) - f(\mathbf{x}) \\ &= \frac{1}{2}(\mathbf{y} - \mathbf{x})^* s^2(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}) \end{aligned}$$

by $t^2 = \|\mathbf{y} - \mathbf{x}\|^2$ and send t to 0, then it follows that

$$0 \leq \frac{1}{2}\mathbf{v}^* d^2 f(\mathbf{x})\mathbf{v}.$$

Because \mathbf{v} is an arbitrary unit vector, the quadratic form $d^2 f(\mathbf{x})$ must be positive semidefinite.

On the other hand, suppose \mathbf{x} is a stationary point of $f(\mathbf{y})$, $d^2 f(\mathbf{x})$ is positive definite, and \mathbf{x} fails to be a local minimum. Then there exists a sequence of points \mathbf{y}_m tending to \mathbf{x} and satisfying

$$0 > f(\mathbf{y}_m) - f(\mathbf{x}) = \frac{1}{2}(\mathbf{y}_m - \mathbf{x})^* s^2(\mathbf{y}_m, \mathbf{x})(\mathbf{y}_m - \mathbf{x}). \quad (5.10)$$

Passing to a subsequence if necessary, we may assume that the unit vectors $\mathbf{v}_m = (\mathbf{y}_m - \mathbf{x})/\|\mathbf{y}_m - \mathbf{x}\|$ converge to a unit vector \mathbf{v} . Dividing inequality (5.10) by $\|\mathbf{y}_m - \mathbf{x}\|^2$ and sending m to ∞ consequently yields $0 \geq \mathbf{v}^* d^2 f(\mathbf{x})\mathbf{v}$, contrary to the hypothesis that $d^2 f(\mathbf{x})$ is positive definite. This contradiction shows that \mathbf{x} represents a local minimum.

To prove the final claim of the proposition, let μ be a nonzero eigenvalue with corresponding eigenvector \mathbf{v} . Then the difference

$$\begin{aligned} f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x}) &= \frac{t^2}{2} \{ \mathbf{v}^* d^2 f(\mathbf{x})\mathbf{v} + \mathbf{v}^* [s^2(\mathbf{x} + t\mathbf{v}, \mathbf{x}) - d^2 f(\mathbf{x})]\mathbf{v} \} \\ &= \frac{t^2}{2} \{ \mu \|\mathbf{v}\|^2 + \mathbf{v}^* [s^2(\mathbf{x} + t\mathbf{v}, \mathbf{x}) - d^2 f(\mathbf{x})]\mathbf{v} \} \end{aligned}$$

has the same sign as μ for t small. ■

Example 5.5.1 Distinguishing Extrema from Saddle Points

Consider the function

$$f(\mathbf{x}) = \frac{1}{4}x_1^4 + \frac{1}{2}x_2^2 - x_1x_2 + x_1 - x_2$$

on \mathbb{R}^2 . It is obvious that

$$\nabla f(\mathbf{x}) = \begin{pmatrix} x_1^3 - x_2 + 1 \\ x_2 - x_1 - 1 \end{pmatrix}, \quad d^2 f(\mathbf{x}) = \begin{pmatrix} 3x_1^2 & -1 \\ -1 & 1 \end{pmatrix}.$$

Adding the two rows of the stationarity equation $\nabla f(\mathbf{x}) = \mathbf{0}$ gives the equation $x_1^3 - x_1 = 0$ with solutions $0, \pm 1$. Solving for x_2 in each case yields

the stationary points $(0, 1)$, $(-1, 0)$, and $(1, 2)$. The last two points are local minima because $d^2f(\mathbf{x})$ is positive definite. The first point is a saddle point because

$$d^2f(0,1) = \begin{pmatrix} 0 & -1 \\ -1 & 1 \end{pmatrix}$$

has characteristic polynomial $\lambda^2 - \lambda - 1$ and eigenvalues $\frac{1}{2}(1 \pm \sqrt{5})$. One of these eigenvalues is positive, and one is negative. ■

We now state and prove a sufficient condition for a point \mathbf{x} to be a constrained local minimum of the objective function $f(\mathbf{y})$. Even in the absence of constraints, inequality (5.11) below represents an improvement over the qualitative claims of Proposition 5.5.1.

Proposition 5.5.2 *Suppose the objective function $f(\mathbf{y})$ of the constrained optimization problem satisfies the multiplier rule (5.1) at the point \mathbf{x} with $\lambda_0 = 1$. Let $f(\mathbf{y})$ and the various constraint functions be twice differentiable at \mathbf{x} , and let $\mathcal{L}(\mathbf{y})$ be the Lagrangian*

$$\mathcal{L}(\mathbf{y}) = f(\mathbf{y}) + \sum_{i=1}^p \lambda_i g_i(\mathbf{y}) + \sum_{j=1}^q \mu_j h_j(\mathbf{y}).$$

If $\mathbf{v}^ d^2\mathcal{L}(\mathbf{x})\mathbf{v} > 0$ for every vector $\mathbf{v} \neq \mathbf{0}$ satisfying $dg_i(\mathbf{x})\mathbf{v} = 0$ and $dh_j(\mathbf{x})\mathbf{v} \leq 0$ for all active constraints, then \mathbf{x} provides a local minimum of $f(\mathbf{y})$. Furthermore, there exists a constant $c > 0$ such that*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + c\|\mathbf{y} - \mathbf{x}\|^2 \quad (5.11)$$

for all feasible \mathbf{y} in a neighborhood of \mathbf{x} .

Proof: Because of the sign restrictions on the μ_j , we have $f(\mathbf{y}) \geq \mathcal{L}(\mathbf{y})$ for all feasible \mathbf{y} in addition to $f(\mathbf{x}) = \mathcal{L}(\mathbf{x})$. If $\mathcal{L}(\mathbf{y}) \geq \mathcal{L}(\mathbf{x})$, then $f(\mathbf{y}) \geq f(\mathbf{x})$, and if $\mathcal{L}(\mathbf{y}) \geq \mathcal{L}(\mathbf{x}) + c\|\mathbf{y} - \mathbf{x}\|^2$, then $f(\mathbf{y}) \geq f(\mathbf{x}) + c\|\mathbf{y} - \mathbf{x}\|^2$. Therefore, it suffices to prove these inequalities for $\mathcal{L}(\mathbf{y})$ rather than $f(\mathbf{y})$. The second inequality $\mathcal{L}(\mathbf{y}) \geq \mathcal{L}(\mathbf{x}) + c\|\mathbf{y} - \mathbf{x}\|^2$ is stronger than the first inequality $\mathcal{L}(\mathbf{y}) \geq \mathcal{L}(\mathbf{x})$, so it also suffices to focus on the second inequality.

With this end in mind, let $s_{\mathcal{L}}^2(\mathbf{y}, \mathbf{x})$ be a second slope function for $\mathcal{L}(\mathbf{y})$ at \mathbf{x} . If the second inequality is false, then there exists a sequence of feasible points \mathbf{y}_m converging to \mathbf{x} and a sequence of positive constants c_m converging to 0 such that

$$\begin{aligned} \mathcal{L}(\mathbf{y}_m) - \mathcal{L}(\mathbf{x}) &= \frac{1}{2}(\mathbf{y}_m - \mathbf{x})^* s_{\mathcal{L}}^2(\mathbf{y}_m, \mathbf{x})(\mathbf{y}_m - \mathbf{x}) \\ &< c_m \|\mathbf{y}_m - \mathbf{x}\|^2. \end{aligned} \quad (5.12)$$

Here $d\mathcal{L}(\mathbf{x})$ vanishes by virtue of the multiplier condition. As usual, we suppose that the sequence of unit vectors

$$\mathbf{v}_m = \frac{1}{\|\mathbf{y}_m - \mathbf{x}\|}(\mathbf{y}_m - \mathbf{x})$$

converges to a unit vector \mathbf{v} by extracting a subsequence if necessary. Dividing inequality (5.12) by $\|\mathbf{y}_m - \mathbf{x}\|^2$ and taking limits then yields $\mathbf{v}^* d^2\mathcal{L}(\mathbf{x})\mathbf{v} \leq 0$. This contradicts our supposition about $d^2\mathcal{L}(\mathbf{x})$ provided we can demonstrate that the tangent conditions $dg_i(\mathbf{x})\mathbf{v} = 0$ and $dh_j(\mathbf{x})\mathbf{v} \leq 0$ hold for all active constraints. These follow by dividing the equations

$$0 = g_i(\mathbf{y}_m) - g_i(\mathbf{x}) = s_{g_i}(\mathbf{y}_m, \mathbf{x})(\mathbf{y}_m - \mathbf{x})$$

and

$$0 \geq h_j(\mathbf{y}_m) - h_j(\mathbf{x}) = s_{h_j}(\mathbf{y}_m, \mathbf{x})(\mathbf{y}_m - \mathbf{x})$$

by $\|\mathbf{y}_m - \mathbf{x}\|$ and taking limits. Recall here that $h_j(\mathbf{x}) = 0$ at an active constraint. ■

Example 5.5.2 *Minimum Eigenvalue of a Symmetric Matrix*

Example 1.4.3 demonstrated how each eigenvector-eigenvalue pair (\mathbf{x}, α) of a symmetric matrix \mathbf{M} provides a stationary point of the Lagrangian

$$\mathcal{L}(\mathbf{x}) = \frac{1}{2}\mathbf{x}^*\mathbf{M}\mathbf{x} - \frac{\alpha}{2}(\|\mathbf{x}\|^2 - 1).$$

Suppose that the eigenvalues are arranged so that $\alpha_1 \leq \dots \leq \alpha_n$ and \mathbf{x}_i is the unit eigenvector corresponding to α_i . We expect that \mathbf{x}_1 furnishes the minimum of $\frac{1}{2}\mathbf{y}^*\mathbf{M}\mathbf{y}$ subject to $g_1(\mathbf{y}) = \frac{1}{2} - \frac{1}{2}\|\mathbf{y}\|^2 = 0$. To check that this is indeed the case, we note that $d^2\mathcal{L}(\mathbf{y}) = \mathbf{M} - \alpha_1\mathbf{I}_n$. The condition $dg_1(\mathbf{x}_1)\mathbf{v} = 0$ is equivalent to $\mathbf{x}_1^*\mathbf{v} = 0$. Because the eigenvectors constitute an orthonormal basis, the equality $\mathbf{x}_1^*\mathbf{v} = 0$ can hold only if

$$\mathbf{v} = \sum_{i=2}^n c_i \mathbf{x}_i.$$

For such a choice of \mathbf{v} , the quadratic form

$$\mathbf{v}^* d^2\mathcal{L}(\mathbf{x}_1)\mathbf{v} = \sum_{i=2}^n c_i^2(\alpha_i - \alpha_1) > 0$$

so long as $\alpha_1 < \alpha_2$ and $\mathbf{v} \neq \mathbf{0}$. Thus, the condition cited in the statement of Proposition 5.5.2 holds, and the point \mathbf{x}_1 minimizes $\mathbf{y}^*\mathbf{M}\mathbf{y}$ subject to the constraint $\|\mathbf{y}\| = 1$.

The case $\alpha_1 = \alpha_2$ is not covered by Proposition 5.5.2. However, we can fall back on an alternative and simpler proof. Consider any unit vector

$$\mathbf{w} = \sum_{i=1}^n b_i \mathbf{x}_i.$$

The inequality

$$\mathbf{w}^* \mathbf{M} \mathbf{w} = \sum_{i=1}^n \alpha_i b_i^2 \|\mathbf{x}_i\|^2 \geq \alpha_1 \sum_{i=1}^n b_i^2 = \alpha_1.$$

immediately establishes the claimed property once we note that $\mathbf{w} = \mathbf{x}_1$ achieves the lower bound α_1 . ■

Example 5.5.3 *Minimum of a Linear Reciprocal Function*

Consider minimizing the nonlinear function $f(\mathbf{x}) = \sum_{i=1}^n c_i x_i^{-1}$ subject to the linear inequality constraint $\sum_{i=1}^n a_i x_i \leq b$. Here all indicated variables and parameters are positive. Differentiating the Lagrangian

$$\mathcal{L}(\mathbf{x}) = \sum_{i=1}^n c_i x_i^{-1} + \mu \left(\sum_{i=1}^n a_i x_i - b \right)$$

gives the multiplier equations

$$-\frac{c_i}{x_i^2} + \mu a_i = 0.$$

It follows that $\mu > 0$, that the constraint is active, and that

$$\begin{aligned} x_i &= \sqrt{\frac{c_i}{\mu a_i}}, & 1 \leq i \leq n \\ \mu &= \left(\frac{1}{b} \sum_{i=1}^n \sqrt{a_i c_i} \right)^2. \end{aligned} \quad (5.13)$$

The second differential $d^2\mathcal{L}(\mathbf{x})$ is diagonal with i th diagonal entry $2c_i/x_i^3$. This matrix is certainly positive definite, and Proposition 5.5.2 confirms that the stationary point (5.13) provides the minimum of $f(\mathbf{x})$ subject to the constraint. ■

When there are only equality constraints, one can say more about the sufficient criterion described in Proposition 5.5.2. Following the discussion in Sect. 5.3, let \mathbf{G} be the $p \times n$ matrix \mathbf{G} with rows $dg_i(\mathbf{x})$ and \mathbf{K} an $n \times (n-p)$ matrix of full rank satisfying $\mathbf{G}\mathbf{K} = \mathbf{0}$. On the kernel of \mathbf{G} the matrix $\mathbf{A} = d^2\mathcal{L}(\mathbf{x})$ is positive definite. Since every \mathbf{v} in the kernel equals some image point $\mathbf{K}\mathbf{u}$, we can establish the validity of the sufficient condition of Proposition 5.5.2 by checking whether the matrix $\mathbf{K}^*\mathbf{A}\mathbf{K}$ of the quadratic form $\mathbf{u}^*\mathbf{K}^*\mathbf{A}\mathbf{K}\mathbf{u}$ is positive definite. There are many practical methods of making this determination. For instance, the sweep operator from computational statistics performs such a check easily in the process of inverting $\mathbf{K}^*\mathbf{A}\mathbf{K}$ [166].

If we want to work directly with the matrix \mathbf{G} , there is another interesting criterion involving the relation between the positive semidefinite matrix

$\mathbf{B} = \mathbf{G}^* \mathbf{G}$ and the second differential $\mathbf{A} = d^2 \mathcal{L}(\mathbf{x})$. One can rephrase the sufficient condition of Proposition 5.5.2 by saying that $\mathbf{v}^* \mathbf{A} \mathbf{v} > 0$ whenever $\mathbf{v}^* \mathbf{B} \mathbf{v} = 0$ and $\mathbf{v} \neq \mathbf{0}$. We claim that this condition is equivalent to the existence of some constant $\gamma > 0$ such that the matrix $\mathbf{A} + \gamma \mathbf{B}$ is positive definite [58]. Clearly, if such a γ exists, then the condition holds. Conversely, suppose the condition holds and that no such γ exists. Then there is a sequence of unit vectors \mathbf{v}_m and a sequence of scalars α_m tending to ∞ such that

$$\mathbf{v}_m^* \mathbf{A} \mathbf{v}_m + \alpha_m \mathbf{v}_m^* \mathbf{B} \mathbf{v}_m \leq 0. \quad (5.14)$$

By passing to a subsequence if needed, we may assume that the sequence \mathbf{v}_m converges to a unit vector \mathbf{v} . On the one hand, because \mathbf{B} is positive semidefinite, inequality (5.14) compels the conclusions $\mathbf{v}_m^* \mathbf{A} \mathbf{v}_m \leq 0$, which must carry over to the limit. On the other hand, dividing inequality (5.14) by α_m and taking limits imply $\mathbf{v}^* \mathbf{B} \mathbf{v} \leq 0$ and therefore $\mathbf{v}^* \mathbf{B} \mathbf{v} = 0$. Because the limit vector \mathbf{v} violates the condition $\mathbf{v}^* \mathbf{A} \mathbf{v} > 0$, the required $\gamma > 0$ exists.

5.6 Problems

1. Find a minimum of $f(\mathbf{x}) = x_1^2 + x_2^2$ subject to the inequality constraints $h_1(\mathbf{x}) = -2x_1 - x_2 + 10 \leq 0$ and $h_2(\mathbf{x}) = -x_1 \leq 0$ on \mathbb{R}^2 . Prove that it is the global minimum.
2. Minimize the function $f(\mathbf{x}) = e^{-(x_1+x_2)}$ subject to the constraints $h_1(\mathbf{x}) = e^{x_1} + e^{x_2} - 20 \leq 0$ and $h_2(\mathbf{x}) = -x_1 \leq 0$ on \mathbb{R}^2 .
3. Find the minimum and maximum of the function $f(\mathbf{x}) = x_1 + x_2$ over the subset of \mathbb{R}^2 defined by the constraints

$$\begin{aligned} h_1(\mathbf{x}) &= -x_1 \\ h_2(\mathbf{x}) &= -x_2 \\ h_3(\mathbf{x}) &= 1 - x_1 x_2. \end{aligned}$$

4. Consider the problem of minimizing $f(\mathbf{x}) = (x_1 + 1)^2 + x_2^2$ subject to the inequality constraint $h(\mathbf{x}) = -x_1^3 + x_2^2 \leq 0$ on \mathbb{R}^2 . Solve the problem by sketching the feasible region and using a little geometry. Show that the multiplier rule of Proposition 5.2.1 with $\lambda_0 = 1$ fails and explain why.
5. In the multiplier rule of Proposition 5.2.1, suppose that the Kuhn-Tucker constraint qualification holds and that $\lambda_0 = 1$. Prove that the remaining multipliers are unique.

6. Consider the inequality constraint functions

$$\begin{aligned}h_1(\mathbf{x}) &= -x_1 \\h_2(\mathbf{x}) &= -x_2 \\h_3(\mathbf{x}) &= x_1^2 + 4x_2^2 - 4 \\h_4(\mathbf{x}) &= (x_1 - 2)^2 + x_2^2 - 5\end{aligned}$$

on \mathbb{R}^2 . Show that the Kuhn-Tucker constraint qualification fails but the Mangasarian-Fromovitz constraint qualification succeeds at the point $x = (0, 1)^*$. For the inequality constraint functions

$$\begin{aligned}h_1(\mathbf{x}) &= x_1^2 - x_2 \\h_2(\mathbf{x}) &= -3x_1^2 + x_2,\end{aligned}$$

show that both constraint qualifications fail at the point $\mathbf{x} = \mathbf{0}$ [96].

7. Consider the two functions $f(\mathbf{x}) = x_1$ and

$$g(\mathbf{x}) = \begin{cases} x_2 & \text{if } x_1 \geq 0 \\ x_2 - x_1^2 & \text{if } x_1 < 0, x_2 \leq 0 \\ x_2 + x_1^2 & \text{if } x_1 < 0, x_2 > 0 \end{cases}$$

defined on \mathbb{R}^2 . Demonstrate that:

- (a) $f(\mathbf{x})$ has differential $df(\mathbf{x}) = (1, 0)$,
- (b) $g(\mathbf{x})$ is continuous except on the half line $\{(x_1, 0)^* : x_1 < 0\}$,
- (c) $g(\mathbf{x})$ has differential $dg(\mathbf{0}) = (0, 1)$ at the origin.

With these functions in hand, it is possible to show that the Lagrange multiplier rule can fail due to lack of continuity of an equality constraint. The optimization problem we have in mind is minimizing $f(\mathbf{x})$ subject to $g(\mathbf{x}) = 0$. Prove that the origin is the unique solution of this problem. In addition prove that no nontrivial pair (μ, λ) satisfies the multiplier condition

$$\mu \nabla f(\mathbf{0}) + \lambda \nabla g(\mathbf{0}) = \mathbf{0}.$$

8. For a real 2×2 matrix

$$\mathbf{M} = \begin{pmatrix} a & b \\ c & d \end{pmatrix},$$

define the Frobenius norm $\|\mathbf{M}\|_F = \sqrt{a^2 + b^2 + c^2 + d^2}$. Let S denote the set of matrices \mathbf{M} with $\det(\mathbf{M}) = 0$. Find the minimum distance from S to the matrix

$$\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix},$$

and exhibit a matrix attaining this distance [69]. (Hints: Introduce a Lagrange multiplier λ in minimizing $\frac{1}{2}\|\mathbf{M}\|_F^2$ subject to $\mathbf{M} \in S$. From the multiplier conditions deduce that $\lambda = \pm 1$ if $b \neq 0$. Show that the assumption $\lambda = \pm 1$ leads to a contradiction. Thus, $b = 0$ and consequently $c = 0$. Express a and d as functions of λ and find the λ 's for which $\det(\mathbf{M}) = 0$.)

9. The equation

$$\sum_{i=1}^n a_i x_i^2 = c$$

defines an ellipse in \mathbb{R}^n whenever all $a_i > 0$. The problem of Apollonius is to find the closest point on the ellipse from an external point \mathbf{y} [30]. Demonstrate that the solution has coordinates

$$x_i = \frac{y_i}{1 + \lambda a_i},$$

where λ is chosen to satisfy

$$\sum_{i=1}^n a_i \left(\frac{y_i}{1 + \lambda a_i} \right)^2 = c.$$

Show how you can adapt this solution to solve the problem with the more general ellipse $(\mathbf{x} - \mathbf{z})^* \mathbf{A} (\mathbf{x} - \mathbf{z}) = c$ for \mathbf{A} a positive definite matrix.

10. Let \mathbf{A} be a positive definite matrix. For a given vector \mathbf{y} , find the maximum of $f(\mathbf{x}) = \mathbf{y}^* \mathbf{x}$ subject to $h(\mathbf{x}) = \mathbf{x}^* \mathbf{A} \mathbf{x} - 1 \leq 0$. Use your result to prove the inequality $|\mathbf{y}^* \mathbf{x}|^2 \leq (\mathbf{x}^* \mathbf{A} \mathbf{x})(\mathbf{y}^* \mathbf{A}^{-1} \mathbf{y})$.
11. Let \mathbf{A} be a full rank $m \times n$ matrix and \mathbf{b} be an $m \times 1$ vector with $m < n$. The set $S = \{\mathbf{x} \in \mathbb{R}^n : \mathbf{A} \mathbf{x} = \mathbf{b}\}$ defines a plane in \mathbb{R}^n . If $m = 1$, S is a hyperplane. Given $\mathbf{y} \in \mathbb{R}^n$, prove that the closest point to \mathbf{y} in S is

$$P(\mathbf{y}) = \mathbf{y} - \mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{A}\mathbf{y} + \mathbf{A}^*(\mathbf{A}\mathbf{A}^*)^{-1}\mathbf{b}.$$

12. If \mathbf{A} is a matrix and \mathbf{y} is a compatible vector, then $\mathbf{A}\mathbf{y} \geq \mathbf{0}$ means that all entries of the vector $\mathbf{A}\mathbf{y}$ are nonnegative. Farkas' lemma says that $\mathbf{x}^* \mathbf{y} \geq 0$ for all vectors \mathbf{y} with $\mathbf{A}\mathbf{y} \geq \mathbf{0}$ if and only if \mathbf{x} is a nonnegative linear combination of the rows of \mathbf{A} . Prove Farkas' lemma assuming that the rows of \mathbf{A} are linearly independent.
13. It is possible to give an elementary proof of the necessity of the Lagrange multiplier rule under a set of alternative conditions [28].

In the setting of Proposition 5.2.1, suppose that the equality constraints $g_i(\mathbf{y})$ are affine and that $f(\mathbf{y})$ and the inequality constraints $h_j(\mathbf{x})$ are merely differentiable at the constrained minimum \mathbf{x} . The constraint qualification is simplified by taking all $h_j(\mathbf{y})$ to be active at \mathbf{x} and requiring the matrix

$$\mathbf{M} = \begin{pmatrix} dg(\mathbf{x}) \\ dh(\mathbf{x}) \end{pmatrix}$$

to have full row rank, where $dg(\mathbf{x})$ and $dh(\mathbf{x})$ stack the differentials of the $g_i(\mathbf{x})$ and $h_j(\mathbf{x})$ row by row. The number of columns of \mathbf{M} should exceed the number of rows of \mathbf{M} . To validate the Lagrange multiplier rule, first argue that the equation

$$\begin{pmatrix} df(\mathbf{x}) \\ \mathbf{M} \end{pmatrix} \mathbf{u} = \begin{pmatrix} c \\ \mathbf{0} \\ \mathbf{v} \end{pmatrix} \quad (5.15)$$

can have no solution \mathbf{u} when c and the entries of \mathbf{v} are all negative numbers. Indeed, if a solution exists, then show that $\mathbf{x} + t\mathbf{u}$ is feasible for $t > 0$ small enough and satisfies $f(\mathbf{x} + t\mathbf{u}) < f(\mathbf{x})$. Next argue that the full rank assumption implies that equation (5.1) holds for some set of Lagrange multipliers. The real question is whether all μ_j are nonnegative. Suppose otherwise and construct a vector \mathbf{v} with all entries negative such that the inner product $\boldsymbol{\mu}^* \mathbf{v}$ is positive. The full rank assumption then implies that the equation

$$\mathbf{M}\mathbf{u} = \begin{pmatrix} \mathbf{0} \\ \mathbf{v} \end{pmatrix}$$

has a solution \mathbf{u} . Finally, demonstrate that taking $c = -\boldsymbol{\mu}^* \mathbf{v}$ forces \mathbf{u} to solve equation (5.15). This contradiction proves that all μ_j are nonnegative. (Hints: Full row rank implies full column rank. Hence, if the matrix appearing on the left side of equation (5.15) has full row rank, then the equation is solvable for any choice of c and \mathbf{v} . This impossibility implies that $df(\mathbf{x})$ can be expressed as a linear combination of the rows of \mathbf{M} .)

14. A random variable takes the value x_i with probability p_i for i ranging from 1 to n . Maximize the entropy $-\sum_{i=1}^n p_i \ln p_i$ subject to a fixed mean $m = \sum_{i=1}^n x_i p_i$. Show that $p_i = \alpha e^{\lambda x_i}$ for constants α and λ . Argue that λ is determined by the equation

$$\sum_{i=1}^n x_i e^{\lambda x_i} = m \sum_{i=1}^n e^{\lambda x_i}.$$

15. Continuing the previous problem, suppose that each $x_j = j$. At the maximum, show that

$$p_i = \frac{p^{i-1}(1-p)}{1-p^n}$$

for some $p > 0$ and all $1 \leq i \leq n$. Argue that p exists and is unique for $n > 1$.

16. Establish the bound

$$\sum_{m=1}^n \left(p_m + \frac{1}{p_m} \right)^2 \geq n^3 + 2n + \frac{1}{n}$$

for a discrete probability density p_1, \dots, p_n . Determine a necessary and sufficient condition for equality to hold [243].

17. Consider the problem of minimizing the continuously differentiable function $f(\mathbf{x})$ subject to the constraint $\mathbf{x} \geq \mathbf{0}$. At a local minimum \mathbf{y} demonstrate that the partial derivative $\partial_i f(\mathbf{y}) = 0$ when $y_i > 0$ and $\partial_i f(\mathbf{y}) \geq 0$ when $y_i = 0$.

18. As a variation on Problem 17, consider minimizing the continuously differentiable function $f(\mathbf{x})$ subject to the constraints $\sum_{i=1}^n x_i = 1$ and $\mathbf{x} \geq \mathbf{0}$. At a local minimum \mathbf{y} demonstrate that there exists a number λ such that the partial derivative $\partial_i f(\mathbf{y}) = \lambda$ when $y_i > 0$ and $\partial_i f(\mathbf{y}) \geq \lambda$ when $y_i = 0$. This result is known as Gibbs' lemma.

19. Prove Nesbitt's inequality

$$\sum_{k=1}^n \frac{p_k}{\sum_{j \neq k} p_j} \geq \frac{n}{n-1}$$

for a discrete probability density p_1, \dots, p_n . Determine a necessary and sufficient condition for equality to hold [243].

20. Find the minimum value of $f(\mathbf{x}) = \|\mathbf{x}\|^2$ subject to the constraints $\sum_{i=1}^n x_i = 1$ and $\mathbf{x} \geq \mathbf{0}$. Interpret the result geometrically.
21. For $p > 1$ define the norm $\|\mathbf{x}\|_p$ on \mathbb{R}^n satisfying $\|\mathbf{x}\|_p^p = \sum_{i=1}^n |x_i|^p$. For a fixed vector \mathbf{z} , maximize $f(\mathbf{x}) = \mathbf{z}^* \mathbf{x}$ subject to $\|\mathbf{x}\|_p^p \leq 1$. Deduce Hölder's inequality $|\mathbf{z}^* \mathbf{x}| \leq \|\mathbf{x}\|_p \|\mathbf{z}\|_q$ for q defined by the equation $p^{-1} + q^{-1} = 1$.
22. Suppose \mathbf{A} is an $n \times n$ positive definite matrix. Find the minimum of $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}$ subject to the constraint $\mathbf{z}^* \mathbf{x} - c \leq 0$. It may help to consider the cases $c \geq 0$ and $c < 0$ separately.

23. Suppose that $\mathbf{v}_1, \dots, \mathbf{v}_m$ are orthogonal eigenvectors of the $n \times n$ symmetric matrix \mathbf{M} . Subject to the constraints

$$\|\mathbf{x}\|^2 = 1, \quad \mathbf{v}_i^* \mathbf{x} = 0, \quad 1 \leq i \leq m < n,$$

show that a minimum of $\mathbf{x}^* \mathbf{M} \mathbf{x}$ must coincide with an eigenvector of \mathbf{M} . Under what circumstances is there a unique minimum of $\mathbf{x}^* \mathbf{M} \mathbf{x}$ subject to the constraints?

24. In the context of Proposition 5.2.1, suppose at the feasible point \mathbf{x} one has $df(\mathbf{x})\mathbf{v} > 0$ for every nontrivial tangent vector \mathbf{v} . Recall that \mathbf{v} satisfies $dg_i(\mathbf{x})\mathbf{v} = 0$ for all equality constraints and $dh_j(\mathbf{x})\mathbf{v} \leq 0$ for all active inequality constraints. Prove that there exists a positive constant c such that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + c\|\mathbf{y} - \mathbf{x}\|$$

for all feasible points \mathbf{y} close enough to \mathbf{x} . Thus, \mathbf{x} represents a local minimum of $f(\mathbf{y})$. (Hints: If the contrary is true, then there exists a sequence of feasible points \mathbf{x}_m such that

$$\|\mathbf{x}_m - \mathbf{x}\| < \frac{1}{m} \text{ and } f(\mathbf{x}_m) < f(\mathbf{x}) + \frac{1}{m}\|\mathbf{x}_m - \mathbf{x}\|.$$

Pass to a subsequence so that $\|\mathbf{x}_m - \mathbf{x}\|^{-1}(\mathbf{x}_m - \mathbf{x})$ converges.)

25. Continuing Problem 24, let $f(\mathbf{x}) = 3x_1 - x_2 + x_1x_2$ and $h_1(\mathbf{x}) = -x_1$, $h_2(\mathbf{x}) = x_1 - x_2$, and $h_3(\mathbf{x}) = x_2 - 2x_1$. Show that $\mathbf{0}$ represents a local minimum by demonstrating the condition $df(\mathbf{0})\mathbf{v} > 0$ for every nontrivial tangent vector \mathbf{v} .
26. In Problem 24, the strict inequality $df(\mathbf{x})\mathbf{v} > 0$ cannot be relaxed to simple inequality. As an example take $f(\mathbf{x}) = x_2$ subject to the constraint $h_1(\mathbf{x}) = -x_1^2 - x_2 \leq 0$. Demonstrate that $df(\mathbf{0})\mathbf{v} \geq 0$ for every vector \mathbf{v} satisfying $dh_1(\mathbf{0})\mathbf{v} \leq 0$, yet $\mathbf{0}$ is not a local minimum.
27. Assume that the objective function $f(\mathbf{x})$ and the equality constraints $g_i(\mathbf{x})$ in an equality constrained minimization problem are continuously differentiable. If the gradients $\nabla g_i(\mathbf{y})$ at a local minimum \mathbf{y} are linearly independent, then the standard Lagrange multiplier rule holds at \mathbf{y} . If in addition $f(\mathbf{x})$ and the $g_i(\mathbf{x})$ possess second differentials at \mathbf{y} , then show that the second differential $d^2\mathcal{L}(\mathbf{y})$ of the Lagrangian satisfies $\mathbf{v}^* d^2\mathcal{L}(\mathbf{x})\mathbf{v} \geq 0$ for every tangent direction \mathbf{v} . (Hints: As demonstrated in Example 4.6.3, there is a one-to-one correspondence between tangent vectors and tangent curves. Expand $\mathcal{L}(\mathbf{x})$ to second order around \mathbf{y} , and use the fact that it coincides with $f(\mathbf{x})$ at feasible points.)

28. Demonstrate that the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ is unbounded below if either (a) \mathbf{A} is not positive semidefinite, or (b) \mathbf{A} is positive semidefinite and $\mathbf{A} \mathbf{x} = -\mathbf{b}$ has no solution. Why does $f(\mathbf{x})$ attain its minimum when it is bounded below? (Hints: Diagonalize \mathbf{A} in the form $\mathbf{O} \mathbf{D} \mathbf{O}^*$, where \mathbf{O} is orthogonal and \mathbf{D} is diagonal. Consider the transformed function $g(\mathbf{z}) = f(\mathbf{x})$ with $\mathbf{z} = \mathbf{O}^* \mathbf{x}$.)
29. Sylvester's criterion states that an $n \times n$ symmetric matrix \mathbf{A} is positive definite if and only if its leading principal minors are positive. To prove this result by induction, verify it in the scalar case $n = 1$. Now consider the $(n + 1) \times (n + 1)$ symmetric block matrix

$$\mathbf{A} = \begin{pmatrix} \mathbf{B} & \mathbf{b} \\ \mathbf{b}^* & c \end{pmatrix},$$

where \mathbf{B} is $n \times n$ and positive definite. Define the function

$$f(\mathbf{x}) = (\mathbf{x}^*, 1) \mathbf{A} \begin{pmatrix} \mathbf{x} \\ 1 \end{pmatrix} = \mathbf{x}^* \mathbf{B} \mathbf{x} + 2\mathbf{b}^* \mathbf{x} + c,$$

and show that $\mathbf{x} = -\mathbf{B}^{-1} \mathbf{b}$ furnishes its minimum. At that point the function has value $c - \mathbf{b}^* \mathbf{B}^{-1} \mathbf{b}$. Thus, $f(\mathbf{x})$ is positive provided $c - \mathbf{b}^* \mathbf{B}^{-1} \mathbf{b}$ is positive. Next verify that

$$\det \mathbf{A} = (c - \mathbf{b}^* \mathbf{B}^{-1} \mathbf{b}) \det \mathbf{B}$$

by checking that

$$\begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{d}^* & 1 \end{pmatrix} \begin{pmatrix} \mathbf{B} & \mathbf{b} \\ \mathbf{b}^* & c \end{pmatrix} \begin{pmatrix} \mathbf{I} & \mathbf{d} \\ \mathbf{0}^* & 1 \end{pmatrix} = \begin{pmatrix} \mathbf{B} & \mathbf{0} \\ \mathbf{0}^* & c - \mathbf{b}^* \mathbf{B}^{-1} \mathbf{b} \end{pmatrix}$$

for an appropriate vector \mathbf{d} . Put all of these hints together, and advance the induction from n to $n + 1$.

30. Let \mathbf{M} be an $m \times n$ matrix. Show that there exist unit vectors \mathbf{u} and \mathbf{v} such that $\mathbf{M} \mathbf{v} = \|\mathbf{M}\| \mathbf{u}$ and $\mathbf{M}^* \mathbf{u} = \|\mathbf{M}\| \mathbf{v}$, where $\|\mathbf{M}\|$ is the matrix norm induced by the Euclidean norms on \mathbb{R}^m and \mathbb{R}^n . (Hint: See Problem 7 in Chap. 2.)
31. Consider the set of $n \times n$ matrices $\mathbf{M} = (m_{ij})$. Demonstrate that $\det \mathbf{M}$ has maximum value $\prod_{i=1}^n d_i$ subject to the constraints

$$\sqrt{\sum_{j=1}^n m_{ij}^2} = d_i$$

for $1 \leq i \leq n$. This is Hadamard's inequality. (Hints: Use the Lagrange multiplier condition and the identities $\det \mathbf{M} = \sum_{j=1}^n m_{ij} M_{ij}$ and $(\mathbf{M}^{-1})_{ij} = M_{ji} / \det \mathbf{M}$ to show that \mathbf{M} can be written as the product $\mathbf{D} \mathbf{R}$ of a diagonal matrix \mathbf{D} with diagonal entries d_i and an orthogonal matrix \mathbf{R} with $\det \mathbf{R} = 1$.)

32. For m a positive integer, verify the explicit second slope

$$s^2(y, x) = 2[y^{m-2} + 2y^{m-3}x + 3y^{m-4}x^2 + \cdots + (m-1)x^{m-2}]$$

of the function $f(x) = x^m$ on \mathbb{R} . Show that

$$\lim_{y \rightarrow x} s^2(y, x) = m(m-1)x^{m-2}.$$

33. Supply the missing algebraic steps in Example 5.4.5.

34. Let $f(y)$ be a real-valued function of the real variable y . Suppose that $f''(y)$ exists at a point x . Prove that

$$f''(x) = \lim_{u \rightarrow 0} \frac{f(x+u) - 2f(x) + f(x-u)}{u^2}.$$

Use Problem 2 of Chap. 4 to devise an example where this limit quotient exists but $f''(x)$ does not exist.

35. Suppose that $f(x)$ is twice differentiable on the interval $(0, \infty)$. If $m_j = \sup_x |f^{(j)}(x)|$, then show that $m_1^2 \leq 4m_0m_2$. (Hints: Expanding $f(x)$ in a second-order Taylor series, demonstrate that

$$|f'(x)| \leq \frac{2m_0}{h} + \frac{m_2h}{2}$$

for all positive h . Choose h to minimize the right-hand side.)

36. Show that the function $f(\mathbf{x}) = e^{x_1} \ln(1+x_2)$ on \mathbb{R}^2 has the second-order Taylor expansion

$$f(\mathbf{x}) = x_2 + x_1x_2 - \frac{1}{2}x_2^2 + R_2(\mathbf{x})$$

around $\mathbf{0}$ with remainder $R_2(\mathbf{x})$.

37. Assume the functions $f(y)$ and $g(y)$ mapping \mathbb{R} into \mathbb{R} are differentiable of order p at the point x . If $f^{(m)}(x) = g^{(m)}(x) = 0$ for all $m < p$ but $g^{(p)}(x) \neq 0$, then demonstrate L'Hôpital's rule

$$\lim_{y \rightarrow x} \frac{f(y)}{g(y)} = \frac{f^{(p)}(x)}{g^{(p)}(x)}.$$

Find the limit of the ratio $\sin^2 x / (e^{x^2} - 1)$ as x tends to 0. (Hint: Consider the Taylor expansion (5.8) for $f(y)$ and the analogous expansion for $g(y)$.)

38. Suppose the function $f(y) : \mathbb{R} \mapsto \mathbb{R}$ is differentiable of order $p > 1$ in a neighborhood of a point x where $f^{(m)}(x) = 0$ for $1 \leq m < p$ and $f^{(p)}(x) \neq 0$. If p is odd, then show that x is a saddlepoint. If p is even, then show that x is a minimum point when $f^{(p)}(x) > 0$ and a maximum point when $f^{(p)}(x) < 0$. (Hint: Invoke the Taylor expansion (5.8).)

6

Convexity

6.1 Introduction

Convexity is one of the cornerstones of mathematical analysis and has interesting consequences for optimization theory, statistical estimation, inequalities, and applied probability. Despite this fact, students seldom see convexity presented in a coherent fashion. It always seems to take a back-seat to more pressing topics. The current chapter is intended as a partial remedy to this pedagogical gap.

We start with convex sets and proceed to convex functions. These intertwined concepts define and illuminate all sorts of inequalities. It is helpful to have a variety of tests to recognize convex functions. We present such tests and discuss the important class of log-convex functions. A strictly convex function has at most one minimum point. This property tremendously simplifies optimization. For a few functions, we are fortunate enough to be able to find their optima explicitly. For other functions, we must iterate.

The definition of a convex function can be extended in various ways. The quasi-convex functions mentioned in the current chapter serve as substitutes for convex functions in many optimization arguments. Later chapters will extend the notion of a convex function to include functions with infinite values. Mathematicians by nature seek to isolate the key properties that drive important theories. However, too much abstraction can be a hindrance in learning. For now we stick to the concrete setting of ordinary convex functions.

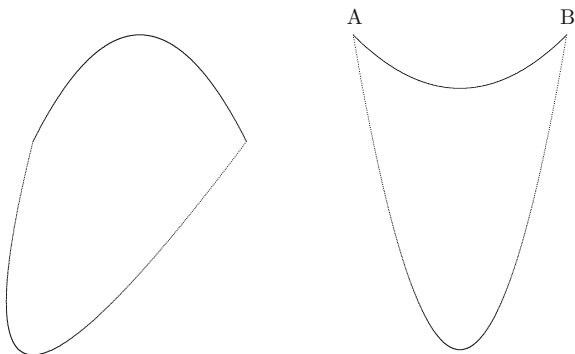


FIGURE 6.1. A convex set on the *left* and a non-convex set on the *right*

The concluding section of this chapter rigorously treats several inequalities from the perspective of probability theory. Our inclusion of Bernstein's proof of the Weierstrass approximation theorem provides a surprising application of Chebyshev's inequality and illustrates the role of probability theory in solving problems outside its usual sphere of influence. The less familiar inequalities of Jensen, Schlömilch, and Hölder find numerous applications in optimization theory and functional analysis.

6.2 Convex Sets

A set $S \subset \mathbb{R}^n$ is said to be convex if the line segment between any two points \mathbf{x} and \mathbf{y} of S lies entirely within S . Formally, this means that whenever $\mathbf{x}, \mathbf{y} \in S$ and $\alpha \in [0, 1]$, the point $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} \in S$ as well. In general, any convex combination $\sum_{i=1}^m \alpha_i \mathbf{x}_i$ of points $\mathbf{x}_1, \dots, \mathbf{x}_m$ from S must also reside in S . Here, the coefficients α_i are nonnegative and sum to 1. Figure 6.1 depicts two sets S in \mathbb{R}^2 , one convex and the other non-convex. The set on the right fails the line segment test for the segment connecting its two cusps A and B.

It is easy to concoct other examples of convex sets. For example, every interval on the real line is convex; every ball in \mathbb{R}^n , either open or closed, is convex; and every multidimensional rectangle, either open, closed, or neither, is convex. Halfspaces and affine subspaces are convex. The former can be open or closed; the latter are always closed. Finitely generated cones as described in Example 2.4.1 are closed convex cones. The set of $n \times n$ positive semidefinite matrices is a closed convex cone in the space of symmetric matrices. The set of $n \times n$ positive definite matrices treated in Example 2.5.4 is an open convex set in the same space. It is not a convex cone because it excludes the $\mathbf{0}$ matrix.

These examples suggest several of the important properties listed in the next proposition.

Proposition 6.2.1 *Convex sets in \mathbb{R}^n enjoy the following properties:*

- (a) *The closure of a convex set is convex.*
- (b) *The interior of a convex set is convex.*
- (c) *A convex set is connected.*
- (d) *The intersection of an arbitrary number of convex sets is convex.*
- (e) *The image and inverse image of a convex set under an affine map $f(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$ are convex.*
- (f) *The Cartesian product $S \times T$ of two convex sets S and T is convex.*

Proof: To prove assertion (a), consider two points \mathbf{x} and \mathbf{y} in the closure of the convex set S . There exist sequences \mathbf{u}_k and \mathbf{v}_k from S converging to \mathbf{x} and \mathbf{y} , respectively. The convex combination $\mathbf{w}_k = \alpha\mathbf{u}_k + (1 - \alpha)\mathbf{v}_k$ is in S as well and converges to $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$ in the closure of S . To verify assertion (b), suppose \mathbf{x} and \mathbf{y} lie in the interior of S . For r sufficiently small, S contains the two balls $\mathbf{x} + B(\mathbf{0}, r)$ and $\mathbf{y} + B(\mathbf{0}, r)$ of radius r . Consider a point $\mathbf{w} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y} + \mathbf{z}$ with $\mathbf{z} \in B(\mathbf{0}, r)$. The decomposition

$$\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} + \mathbf{z} = \alpha(\mathbf{x} + \mathbf{z}) + (1 - \alpha)(\mathbf{y} + \mathbf{z})$$

makes it obvious that \mathbf{w} also lies in S . Assertion (c) is a consequence of the fact that an arcwise connected set is connected. Assertions (d), (e), and (f) follow directly from the definitions. ■

Some obvious corollaries can be drawn from items (e) and (f) of Proposition 6.2.1. For example, suppose S and T are convex sets and λ is any real number. Then we can assert that the sets λS and $S + T$ are convex. Convex sets have many other crucial properties. Here is one that figures prominently in optimization theory.

Proposition 6.2.2 *For a convex set S of \mathbb{R}^n , there is at most one point $\mathbf{y} \in S$ attaining the minimum distance $\text{dist}(\mathbf{x}, S)$ from \mathbf{x} to S . If S is closed, there is exactly one point.*

Proof: These claims are obvious if \mathbf{x} is in S . Suppose that \mathbf{x} is not in S and that \mathbf{y} and \mathbf{z} in S both attain the minimum. Then $\frac{1}{2}(\mathbf{y} + \mathbf{z}) \in S$ and

$$\begin{aligned} \text{dist}(\mathbf{x}, S) &\leq \left\| \mathbf{x} - \frac{1}{2}(\mathbf{y} + \mathbf{z}) \right\| \\ &\leq \frac{1}{2}\|\mathbf{x} - \mathbf{y}\| + \frac{1}{2}\|\mathbf{x} - \mathbf{z}\| \\ &= \text{dist}(\mathbf{x}, S). \end{aligned}$$

Hence, equality must hold in the displayed triangle inequality. This is possible if and only if $\mathbf{x} - \mathbf{y} = c(\mathbf{x} - \mathbf{z})$ for some positive number c . In view of the fact that $\text{dist}(\mathbf{x}, S) = \|\mathbf{x} - \mathbf{y}\| = \|\mathbf{x} - \mathbf{z}\|$, the value of c is 1, and \mathbf{y} and \mathbf{z} coincide. The second claim follows from Example 2.5.5 of Chap. 2. ■

Another important property relates to separation by hyperplanes.

Proposition 6.2.3 *Consider a closed convex set S of \mathbb{R}^n and a point \mathbf{x} outside S . There exists a unit vector \mathbf{v} and real number c such that*

$$\mathbf{v}^* \mathbf{x} > c \geq \mathbf{v}^* \mathbf{z} \quad (6.1)$$

for all $\mathbf{z} \in S$. As a consequence, S equals the intersection of all closed halfspaces containing it. If \mathbf{x} is a boundary point of a convex set S , then there exists a unit vector \mathbf{v} such that $\mathbf{v}^* \mathbf{x} \geq \mathbf{v}^* \mathbf{z}$ for all $\mathbf{z} \in S$.

Proof: Let \mathbf{y} be the closest point to \mathbf{x} in S . Suppose that we can prove the obtuse angle criterion

$$(\mathbf{x} - \mathbf{y})^*(\mathbf{z} - \mathbf{y}) \leq 0 \quad (6.2)$$

for all $\mathbf{z} \in S$. If we take $\mathbf{v} = \mathbf{x} - \mathbf{y}$, then $c = \mathbf{v}^* \mathbf{y} \geq \mathbf{v}^* \mathbf{z}$ for all $\mathbf{z} \in S$. Furthermore, $\mathbf{v}^* \mathbf{x} > \mathbf{v}^* \mathbf{y} = c$ because $\mathbf{v}^* \mathbf{v} = \|\mathbf{v}\|^2 > 0$. One can clearly replace \mathbf{v} by $\|\mathbf{v}\|^{-1} \mathbf{v}$ without disrupting the separation inequalities (6.1).

To prove inequality (6.2), suppose it fails. Then $(\mathbf{x} - \mathbf{y})^*(\mathbf{z} - \mathbf{y}) > 0$ for some $\mathbf{z} \in S$. For each $0 < \alpha < 1$, the point $\alpha \mathbf{z} + (1 - \alpha)\mathbf{y}$ is in S , and

$$\begin{aligned} \|\mathbf{x} - \alpha \mathbf{z} - (1 - \alpha)\mathbf{y}\|^2 &= \|\mathbf{x} - \mathbf{y} - \alpha(\mathbf{z} - \mathbf{y})\|^2 \\ &= \|\mathbf{x} - \mathbf{y}\|^2 - \alpha [2(\mathbf{x} - \mathbf{y})^*(\mathbf{z} - \mathbf{y}) - \alpha \|\mathbf{z} - \mathbf{y}\|^2]. \end{aligned}$$

For α sufficiently small, the term above in square brackets is positive, so $\alpha \mathbf{z} + (1 - \alpha)\mathbf{y}$ improves on the choice of \mathbf{y} . This contradiction demonstrates inequality (6.2).

If \mathbf{x} is a boundary point of S , then there exists a sequence of points \mathbf{x}_i outside S that converge to \mathbf{x} . Let \mathbf{v}_i be the unit vector defining the hyperplane separating \mathbf{x}_i from S . Without loss of generality, we can assume that some subsequence \mathbf{v}_{i_j} converges to a unit vector \mathbf{v} . Taking limits in the strict inequality $\mathbf{v}_{i_j}^* \mathbf{x}_{i_j} > \mathbf{v}_{i_j}^* \mathbf{z}$ for $\mathbf{z} \in S$ then yields the desired result $\mathbf{v}^* \mathbf{x} \geq \mathbf{v}^* \mathbf{z}$. ■

Example 6.2.1 *Farkas' Lemma and Markov Chains*

As a continuation of Example 2.4.1, consider the finitely generated cone

$$C = \left\{ \sum_{i=1}^m \alpha_i \mathbf{v}_i : \alpha_i \geq 0, i = 1, \dots, m \right\}.$$

Because C is closed and convex, any point \mathbf{x} outside C can be separated from C by a hyperplane. Thus, there exists a vector \mathbf{w} and a constant

b with $\mathbf{w}^* \mathbf{x} > b \geq \mathbf{w}^* \mathbf{y}$ for all \mathbf{y} in C . Given the origin $\mathbf{0}$ belongs to C , the constant $b \geq 0$. In fact, b must equal 0. If $\mathbf{w}^* \mathbf{y} > 0$ for some \mathbf{y} in C , then $r\mathbf{w}^* \mathbf{y} > b$ for some $r > 0$. Since $r\mathbf{y}$ belongs to C , we reach a contradiction. Farkas' lemma summarizes these findings by posing two mutually exclusive alternatives, one of which must hold. Either the point \mathbf{x} satisfies $\mathbf{x} = \sum_{i=1}^m \alpha_i \mathbf{v}_i$ for nonnegative constants α_i , or there exists a vector \mathbf{w} with $\mathbf{w}^* \mathbf{x} > 0$ and $0 \geq \mathbf{w}^* \mathbf{v}_i$ for all i .

Farkas' lemma has a clever application in Markov chain theory [97]. Recall that a Markov chain on n states is governed by a transition probability matrix $\mathbf{P} = (p_{ij})$ with nonnegative entries and row sums equal to 1. A stationary distribution is a row vector $\boldsymbol{\alpha}$ with nonnegative entries α_i that sum equal to 1 and satisfy the equilibrium condition $\boldsymbol{\alpha} = \boldsymbol{\alpha} \mathbf{P}$. Farkas' lemma gives an easy proof that such vectors $\boldsymbol{\alpha}$ exist. The stationarity conditions can be restated as the vector equation

$$\begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} = \sum_{i=1}^n \alpha_i \mathbf{v}_i$$

for vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ with entries $v_{ij} = p_{ij} - 1_{\{j=i\}}$ and $v_{i,n+1} = 1$ for all i and j between 1 and n . The Farkas alternative postulates the existence of a vector $\mathbf{w} = (w_1, \dots, w_{n+1})^*$ with

$$\begin{aligned} \mathbf{w}^* \begin{pmatrix} \mathbf{0} \\ 1 \end{pmatrix} &= w_{n+1} > 0 \\ \mathbf{w}^* \mathbf{v}_i &= \sum_{j=1}^n w_j p_{ij} - w_i + w_{n+1} \leq 0 \end{aligned}$$

for $1 \leq i \leq n$. Choosing $w_i = \min_{1 \leq j \leq n} w_j$, we find that

$$0 \geq \sum_{j=1}^n w_j p_{ij} - w_i + w_{n+1} \geq w_{n+1},$$

contradicting the assumption $w_{n+1} > 0$. Thus, the stationarity conditions must hold. \blacksquare

We now turn to a famous convexity result of Carathéodory and its application to compact sets. The convex hull of a set S , denoted $\text{conv } S$, is the smallest convex set containing S . Equivalently, $\text{conv } S$ consists of all convex combinations $\sum_i \alpha_i \mathbf{v}_i$ of points \mathbf{v}_i from S . The convex hull displayed in Fig. 6.2 consists of the boundary plus all points internal to it.

Proposition 6.2.4 (Carathéodory) *Let S be a nonempty subset of \mathbb{R}^n . Every vector from $\text{conv } S$ can be represented as a convex combination of at most $n + 1$ vectors from S . Furthermore, if S is compact, then $\text{conv } S$ is also compact.*

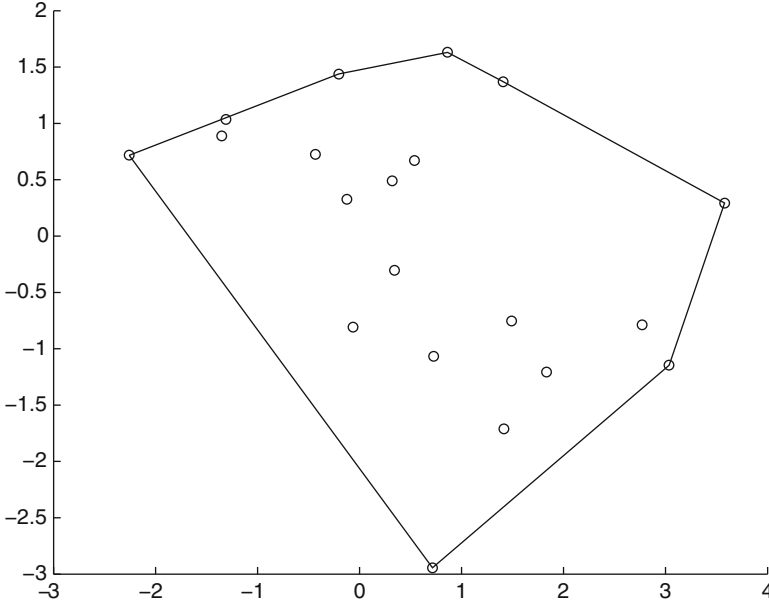


FIGURE 6.2. Convex hull of 20 random points in the plane

Proof: Consider the set $T = \{(\mathbf{y}, 1) : \mathbf{y} \in S\}$. A point $(\mathbf{x}, 1)$ in $\text{conv} T$ can be represented as a convex combination $(\mathbf{x}, 1) = \sum_i \alpha_i (\mathbf{v}_i, 1)$. The point $(\mathbf{x}, 1)$ also belongs to the cone generated by the vectors $(\mathbf{v}_i, 1)$. As noted in Example 2.4.1, we can eliminate all but $n + 1$ linearly independent vectors $(\mathbf{v}_{i_j}, 1)$ in this representation. It follows that $\mathbf{x} = \sum_j \beta_j \mathbf{v}_{i_j}$ with $1 = \sum_j \beta_j$ and all $\beta_j \geq 0$.

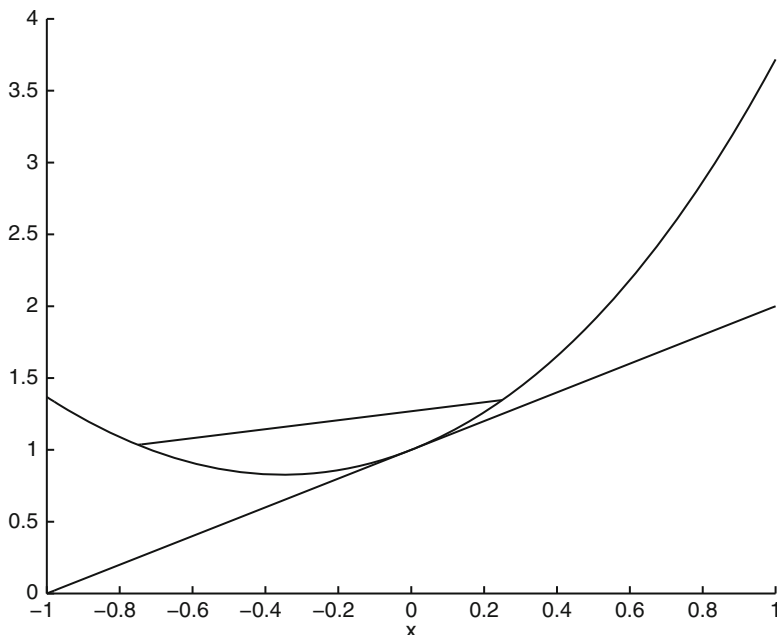
When S is compact, consider a sequence $\mathbf{x}_j = \sum_i \beta_{ji} \mathbf{v}_{ji}$ from $\text{conv} S$. By the first part of the proposition, one can assume that the sum over i runs from 1 to $n + 1$ at most. It suffices to prove that \mathbf{x}_j possesses a subsequence that converges to a point in $\text{conv} S$. By passing to successive subsequences as needed, one can assume that each β_{ji} converges to $\beta_i \geq 0$ and each \mathbf{v}_{ji} converges to $\mathbf{v}_i \in S$. It follows that \mathbf{x}_j converges to the convex combination $\sum_i \beta_i \mathbf{v}_i$ of points from S . ■

6.3 Convex Functions

Convex functions are defined on convex sets. A real-valued function $f(\mathbf{x})$ defined on a convex set S is convex provided

$$f[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] \leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) \tag{6.3}$$

for all $\mathbf{x}, \mathbf{y} \in S$ and $\alpha \in [0, 1]$. Figure 6.3 depicts how in one dimension definition (6.3) requires the chord connecting two points on the curve $f(x)$

FIGURE 6.3. Plot of the convex function $e^x + x^2$

to lie above the curve. If strict inequality holds in inequality (6.3) for every $\mathbf{x} \neq \mathbf{y}$ and $\alpha \in (0, 1)$, then $f(\mathbf{x})$ is said to be strictly convex. One can prove by induction that inequality (6.3) extends to

$$f\left(\sum_{i=1}^m \alpha_i \mathbf{x}_i\right) \leq \sum_{i=1}^m \alpha_i f(\mathbf{x}_i)$$

for any convex combination of points from S . This is the finite form of Jensen's inequality. Proposition 6.6.1 discusses an integral form. A concave function satisfies the reverse of inequality (6.3).

Example 6.3.1 *Affine Functions Are Convex*

For an affine function $f(\mathbf{x}) = \mathbf{a}^* \mathbf{x} + b$, equality holds in inequality (6.3). ■

Example 6.3.2 *Norms Are Convex*

The Euclidean norm $f(\mathbf{x}) = \|\mathbf{x}\| = \sqrt{\sum_{i=1}^n x_i^2}$ satisfies the triangle inequality and the homogeneity condition $\|c\mathbf{x}\| = |c| \|\mathbf{x}\|$. Thus,

$$\|\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}\| \leq \|\alpha \mathbf{x}\| + \|(1 - \alpha) \mathbf{y}\| = \alpha \|\mathbf{x}\| + (1 - \alpha) \|\mathbf{y}\|$$

for every $\alpha \in [0, 1]$. The same argument works for any norm. The choice $\mathbf{y} = 2\mathbf{x}$ gives equality in inequality (6.3) and shows that no norm is strictly convex. ■

Example 6.3.3 *The Distance to a Convex Set Is Convex*

The distance $\text{dist}(\mathbf{x}, S)$ from a point $\mathbf{x} \in \mathbb{R}^n$ to a convex set S is convex in \mathbf{x} . Indeed, for any convex combination $\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$, take sequences \mathbf{u}_k and \mathbf{v}_k from S such that

$$\begin{aligned}\text{dist}(\mathbf{x}, S) &= \lim_{k \rightarrow \infty} \|\mathbf{x} - \mathbf{u}_k\| \\ \text{dist}(\mathbf{y}, S) &= \lim_{k \rightarrow \infty} \|\mathbf{y} - \mathbf{v}_k\|.\end{aligned}$$

The points $\alpha\mathbf{u}_k + (1 - \alpha)\mathbf{v}_k$ lie in S , and taking limits in the inequality

$$\begin{aligned}\text{dist}[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}, S] &\leq \|\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} - \alpha\mathbf{u}_k - (1 - \alpha)\mathbf{v}_k\| \\ &\leq \alpha\|\mathbf{x} - \mathbf{u}_k\| + (1 - \alpha)\|\mathbf{y} - \mathbf{v}_k\|\end{aligned}$$

yields $\text{dist}[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}, S] \leq \alpha \text{dist}(\mathbf{x}, S) + (1 - \alpha) \text{dist}(\mathbf{y}, S)$. ■

Example 6.3.4 *Convex Functions Generate Convex Sets*

Consider a convex function $f(\mathbf{x})$ defined on \mathbb{R}^n . Examination of definition (6.3) shows that the sublevel sets $\{\mathbf{x} : f(\mathbf{x}) \leq c\}$ and $\{\mathbf{x} : f(\mathbf{x}) < c\}$ are convex for any constant c . They may be empty. Conversely, a closed convex set S can be represented as $\{\mathbf{x} : f(\mathbf{x}) \leq 0\}$ using the continuous convex function $f(\mathbf{x}) = \text{dist}(\mathbf{x}, S)$. This result does not preclude the possibility that a convex set is a sublevel set of a nonconvex function. For instance, the set $\{\mathbf{x} : 1 - x_1x_2 \leq 0, x_1 \geq 0, x_2 \geq 0\}$ is convex while the function $1 - x_1x_2$ is nonconvex on the domain $\{\mathbf{x} : x_1 \geq 0, x_2 \geq 0\}$. ■

Example 6.3.5 *A Convex Function Has a Convex Epigraph*

The epigraph of a real-valued function $f(\mathbf{x})$ is defined as the set of points (\mathbf{y}, r) with $f(\mathbf{y}) \leq r$. Roughly speaking, the epigraph is the region lying above the graph of $f(\mathbf{x})$. Consider two points (\mathbf{y}, r) and (\mathbf{z}, s) in the epigraph of $f(\mathbf{x})$. If $f(\mathbf{x})$ is convex, then

$$\begin{aligned}f[\alpha\mathbf{y} + (1 - \alpha)\mathbf{z}] &\leq \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{z}) \\ &\leq \alpha r + (1 - \alpha)s,\end{aligned}$$

and the convex combination $\alpha(\mathbf{y}, r) + (1 - \alpha)(\mathbf{z}, s)$ occurs in the epigraph of $f(\mathbf{x})$. Conversely, if the epigraph of $f(\mathbf{x})$ is a convex set, then $f(\mathbf{x})$ must be a convex function. ■

Figure 6.3 illustrates how a tangent line to a convex curve lies below the curve. This property characterizes convex differentiable functions.

Proposition 6.3.1 *Let $f(\mathbf{x})$ be a differentiable function on the open convex set $S \subset \mathbb{R}^n$. Then $f(\mathbf{x})$ is convex if and only if*

$$f(\mathbf{y}) \geq f(\mathbf{x}) + df(\mathbf{x})(\mathbf{y} - \mathbf{x}) \tag{6.4}$$

for all $\mathbf{x}, \mathbf{y} \in S$. Furthermore, $f(\mathbf{x})$ is strictly convex if and only if strict inequality prevails in inequality (6.4) when $\mathbf{y} \neq \mathbf{x}$.

Proof: If $f(\mathbf{x})$ is convex, then we can rearrange inequality (6.3) to give

$$\begin{aligned} \frac{f[\mathbf{x} + (1 - \alpha)(\mathbf{y} - \mathbf{x})] - f(\mathbf{x})}{(1 - \alpha)} &= \frac{f[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}] - f(\mathbf{x})}{1 - \alpha} \\ &\leq f(\mathbf{y}) - f(\mathbf{x}). \end{aligned}$$

Letting α tend to 1 proves inequality (6.4). To demonstrate the converse, let $\mathbf{z} = \alpha\mathbf{x} + (1 - \alpha)\mathbf{y}$. Then with obvious notational changes, inequality (6.4) implies

$$\begin{aligned} f(\mathbf{x}) &\geq f(\mathbf{z}) + df(\mathbf{z})(\mathbf{x} - \mathbf{z}) \\ f(\mathbf{y}) &\geq f(\mathbf{z}) + df(\mathbf{z})(\mathbf{y} - \mathbf{z}). \end{aligned}$$

Multiplying the first of these inequalities by α and the second by $1 - \alpha$ and adding the results produce

$$\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) \geq f(\mathbf{z}) + df(\mathbf{z})(\mathbf{z} - \mathbf{z}) = f(\mathbf{z}),$$

which is just inequality (6.3). The claims about strict convexity are left to the reader. ■

It is useful to have simpler tests for convexity than inequalities (6.3) and (6.4). One such test involves the second differential $d^2f(\mathbf{x})$ of a function $f(\mathbf{x})$.

Proposition 6.3.2 *Consider a twice differentiable function $f(\mathbf{x})$ on the open convex set $S \subset \mathbb{R}^n$. If its second differential $d^2f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} , then $f(\mathbf{x})$ is convex. If $d^2f(\mathbf{x})$ is positive definite for all \mathbf{x} , then $f(\mathbf{x})$ is strictly convex.*

Proof: The expansion

$$\begin{aligned} f(\mathbf{y}) &= f(\mathbf{x}) + df(\mathbf{x})(\mathbf{y} - \mathbf{x}) \\ &\quad + (\mathbf{y} - \mathbf{x})^* \int_0^1 d^2f[\mathbf{x} + t(\mathbf{y} - \mathbf{x})](1 - t) dt (\mathbf{y} - \mathbf{x}) \end{aligned}$$

for $\mathbf{y} \neq \mathbf{x}$ shows that

$$f(\mathbf{y}) \geq f(\mathbf{x}) + df(\mathbf{x})(\mathbf{y} - \mathbf{x}),$$

with strict inequality when $d^2f(\mathbf{x})$ is positive definite for all \mathbf{x} . ■

Example 6.3.6 *Generalized Arithmetic-Geometric Mean Inequality*

The second derivative test shows that the function e^x is strictly convex. Taking $y_i = e^{x_i}$, $\sum_{i=1}^n \alpha_i = 1$, and all $\alpha_i \geq 0$ produces the generalized arithmetic-geometric mean inequality

$$\prod_{i=1}^n y_i^{\alpha_i} \leq \sum_{i=1}^n \alpha_i y_i. \quad (6.5)$$

Equality holds if all y_i coincide or all but one α_i equals 0. ■

Example 6.3.7 *Strictly Convex Quadratic Functions*

If the matrix \mathbf{A} is positive definite, then Proposition 6.3.2 implies that the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^*\mathbf{A}\mathbf{x} + \mathbf{b}^*\mathbf{x} + c$ is strictly convex. ■

Even Proposition 6.3.2 can be difficult to apply. The next proposition helps us to recognize convex functions by their closure properties.

Proposition 6.3.3 *Convex functions satisfy the following:*

- (a) *If $f(\mathbf{x})$ is convex and $g(\mathbf{x})$ is convex and increasing, then the functional composition $g \circ f(\mathbf{x})$ is convex.*
- (b) *If $f(\mathbf{x})$ is convex, then the functional composition $f(\mathbf{A}\mathbf{x} + \mathbf{b})$ of $f(\mathbf{x})$ with an affine function $\mathbf{A}\mathbf{x} + \mathbf{b}$ is convex.*
- (c) *If $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex and α and β are nonnegative constants, then $\alpha f(\mathbf{x}) + \beta g(\mathbf{x})$ is convex.*
- (d) *If $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex, then $\max\{f(\mathbf{x}), g(\mathbf{x})\}$ is convex.*
- (e) *If $f_m(\mathbf{x})$ is a sequence of convex functions, then $\lim_{m \rightarrow \infty} f_m(\mathbf{x})$ is convex whenever it exists.*

Proof: To prove assertion (a), we calculate

$$\begin{aligned} g \circ f[\alpha\mathbf{x} + (1 - \alpha)\mathbf{y}] &\leq g[\alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y})] \\ &\leq \alpha g \circ f(\mathbf{x}) + (1 - \alpha)g \circ f(\mathbf{y}). \end{aligned}$$

The remaining assertions are left to the reader. ■

Part (a) of Proposition 6.3.3 implies that $e^{f(\mathbf{x})}$ is convex when $f(\mathbf{x})$ is convex and that $f(\mathbf{x})^\beta$ is convex when $f(\mathbf{x})$ is nonnegative and convex and $\beta > 1$. One case not covered by the Proposition is products. The counterexample $x^3 = x^2x$ shows that the product of two convex functions is not necessarily convex. In some situations the limit of a sequence of convex functions is no longer finite. Many authors consider $+\infty$ to be a legitimate value for a convex function while $-\infty$ is illegitimate. For the sake of simplicity, we prefer to deal with functions having only finite values. In Chap. 14 we relax this restriction.

Example 6.3.8 *Largest Eigenvalue of a Symmetric Matrix*

Part (d) of Proposition 6.3.3 can be generalized. Suppose the function $f(\mathbf{x}, \mathbf{y})$ is convex in \mathbf{y} for each fixed \mathbf{x} . Then, provided it is finite, the function $\sup_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$ is convex in \mathbf{y} . For a specific application, recall that Example 1.4.3 proves the formula

$$\lambda_{\max}(\mathbf{M}) = \max_{\|\mathbf{x}\|=1} \mathbf{x}^*\mathbf{M}\mathbf{x}$$

for the largest eigenvalue of a symmetric matrix. Because the map taking \mathbf{M} into $\mathbf{x}^* \mathbf{M} \mathbf{x}$ is linear, it follows that $\lambda_{\max}(\mathbf{M})$ is convex in \mathbf{M} . Applying the same reasoning to $-\mathbf{M}$, we deduce that the minimum eigenvalue $\lambda_{\min}(\mathbf{M})$ is concave in \mathbf{M} . ■

Example 6.3.9 Differences of Convex Functions

Although the class of convex functions is rather narrow, most well-behaved functions can be expressed as the difference of two convex functions. For example, consider a polynomial $p(x) = \sum_{m=0}^n p_m x^m$. The second derivative test shows that x^m is convex whenever m is even. If m is odd, then x^m is convex on $[0, \infty)$, and $-x^m$ is convex on $(-\infty, 0)$. Therefore,

$$x^m = \max\{x^m, 0\} - \max\{-x^m, 0\}$$

is the difference of two convex functions. Because the class of differences of convex functions is closed under the formation of linear combinations, it follows that $p(x)$ belongs to this larger class. ■

A positive function $f(\mathbf{x})$ is said to be log-convex if $\ln f(\mathbf{x})$ is convex. Log-convex functions have excellent closure properties as documented by the next proposition.

Proposition 6.3.4 *Log-convex functions satisfy the following:*

- (a) *If $f(\mathbf{x})$ is log-convex, then $f(\mathbf{x})$ is convex.*
- (b) *If $f(\mathbf{x})$ is convex and $g(\mathbf{x})$ is log-convex and increasing, then the functional composition $g \circ f(\mathbf{x})$ is log-convex.*
- (c) *If $f(\mathbf{x})$ is log-convex, then the functional composition $f(\mathbf{A}\mathbf{x} + \mathbf{b})$ of $f(\mathbf{x})$ with an affine function $\mathbf{A}\mathbf{x} + \mathbf{b}$ is log-convex.*
- (d) *If $f(\mathbf{x})$ is log-convex, then $f(\mathbf{x})^\alpha$ and $\alpha f(\mathbf{x})$ are log-convex for any $\alpha > 0$.*
- (e) *If $f(\mathbf{x})$ and $g(\mathbf{x})$ are log-convex, then $f(\mathbf{x}) + g(\mathbf{x})$, $\max\{f(\mathbf{x}), g(\mathbf{x})\}$, and $f(\mathbf{x})g(\mathbf{x})$ are log-convex.*
- (f) *If $f_m(\mathbf{x})$ is a sequence of log-convex functions, then $\lim_{m \rightarrow \infty} f_m(\mathbf{x})$ is log-convex whenever it exists and is positive.*

Proof: Assertion (a) follows from part (a) of Proposition 6.3.3 after composing the functions e^x and $\ln f(\mathbf{x})$. To prove that the sum of log-convex functions is log-convex, we let $h(\mathbf{x}) = f(\mathbf{x}) + g(\mathbf{x})$ and apply Hölder's inequality as stated in Problem 21 of Chap. 5 and in Example 6.6.3 later in this chapter. Taking $\alpha = 1/p$ and $1 - \alpha = 1/q$ consequently implies that

$$\begin{aligned}
 h[\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}] &= f[\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}] + g[\alpha \mathbf{x} + (1 - \alpha) \mathbf{y}] \\
 &\leq f(\mathbf{x})^\alpha f(\mathbf{y})^{1-\alpha} + g(\mathbf{x})^\alpha g(\mathbf{y})^{1-\alpha} \\
 &\leq [f(\mathbf{x}) + g(\mathbf{x})]^\alpha [f(\mathbf{y}) + g(\mathbf{y})]^{1-\alpha} \\
 &= h(\mathbf{x})^\alpha h(\mathbf{y})^{1-\alpha}.
 \end{aligned}$$

The remaining assertions are left to the reader. ■

Example 6.3.10 *The Convex Function of Gordon’s Theorem*

In Proposition 5.3.2, we encountered the function

$$f(\mathbf{x}) = \ln \left[\sum_{j=1}^r \exp(\mathbf{z}_j^* \mathbf{x}) \right].$$

Given the log-convexity of the functions $\exp(\mathbf{z}_j^* \mathbf{x})$, we now recognize $f(\mathbf{x})$ as convex. This is one of the reasons for its success in Gordon’s theorem. ■

Example 6.3.11 *Gamma Function*

Gauss’s representation of the gamma function

$$\Gamma(z) = \lim_{n \rightarrow \infty} \frac{n! n^z}{z(z+1) \cdots (z+n)} \tag{6.6}$$

shows that it is log-convex on $(0, \infty)$ [132]. Indeed, one can easily check that n^z and $(z+k)^{-1}$ are log-convex and then apply the closure of the set of log-convex functions under the formation of products and limits. Note that invoking convexity in this argument is insufficient because the set of convex functions is not closed under the formation of products. Alternatively, one can deduce log-convexity from Euler’s definition

$$\Gamma(z) = \int_0^\infty x^{z-1} e^{-x} dx$$

by viewing the integral as the limit of Riemann sums, each of which is log-convex. ■

Example 6.3.12 *Log-concavity of $\det \Sigma$ for Σ Positive Definite*

Let Ω be an $n \times n$ positive definite matrix. According to Appendix A.2, the function

$$f(\mathbf{y}) = \left(\frac{1}{2\pi} \right)^{n/2} |\det \Omega|^{-1/2} e^{-\mathbf{y}^* \Omega^{-1} \mathbf{y}/2}$$

is a probability density. Integrating over all $\mathbf{y} \in \mathbb{R}^n$ produces the identity

$$|\det \Omega|^{1/2} = \frac{1}{(2\pi)^{n/2}} \int e^{-\mathbf{y}^* \Omega^{-1} \mathbf{y}/2} d\mathbf{y}.$$

We can restate this identity in terms of the inverse matrix $\Sigma = \Omega^{-1}$ as

$$\ln \det \Sigma = n \ln(2\pi) - 2 \ln \int e^{-\mathbf{y}^* \Sigma \mathbf{y} / 2} d\mathbf{y}.$$

By the reasoning of the last two examples, the integral on the right is log-convex. Because Σ is positive definite if and only if Ω is positive definite, it follows that $\ln \det \Sigma$ is concave in the positive definite matrix Σ . ■

6.4 Continuity, Differentiability, and Integrability

In this section we discuss some continuity, differentiability, and integrability properties of convex functions. Let us start with a sobering counterexample involving the closed unit ball $C(\mathbf{0}, 1)$ of \mathbb{R}^n and a positive function $g(\mathbf{y})$ defined on its boundary. One can extend $g(\mathbf{y})$ to a convex function $f(\mathbf{x})$ with domain $C(\mathbf{0}, 1)$ by setting

$$f(\mathbf{x}) = \begin{cases} 0 & \|\mathbf{x}\| < 1 \\ g(\mathbf{x}) & \|\mathbf{x}\| = 1. \end{cases}$$

Even though $f(\mathbf{x})$ is convex throughout $C(\mathbf{0}, 1)$, it is discontinuous everywhere on the boundary. Even worse, $f(\mathbf{x})$ is not even lower semicontinuous. As the next proposition demonstrates, matters improve considerably if we restrict our attention to the interior of the domain of a convex function.

Proposition 6.4.1 *A convex function $f(\mathbf{x})$ is continuous on the interior of its domain and locally Lipschitz around every interior point. In other words, there exists a constant c such that $|f(\mathbf{z}) - f(\mathbf{y})| \leq c\|\mathbf{z} - \mathbf{y}\|$ for all \mathbf{y} and \mathbf{z} near \mathbf{x} .*

Proof: Let \mathbf{y} be an interior point and $C(\mathbf{y}, r)$ be a closed ball of radius r around \mathbf{y} contained within the domain of $f(\mathbf{x})$. Without loss of generality, we may assume that $\mathbf{y} = \mathbf{0}$. We first demonstrate that $f(\mathbf{x})$ is bounded above near $\mathbf{0}$. Define the $n + 1$ points

$$\mathbf{v}_0 = -\frac{r}{2n}\mathbf{1}, \quad \mathbf{v}_i = r\mathbf{e}_i - \frac{r}{2n}\mathbf{1}, \quad 1 \leq i \leq n,$$

using the standard basis \mathbf{e}_i of \mathbb{R}^n . It is easy to check that all of these points lie in $C(\mathbf{0}, r)$. Hence, any convex combination $\sum_{i=0}^n \alpha_i \mathbf{v}_i$ also lies in $C(\mathbf{0}, r)$. Even more surprising, any point \mathbf{x} in the open interval

$$J = \prod_{i=1}^n \left(-\frac{r}{2n}, \frac{r}{2n}\right)$$

can be represented as such a convex combination. This assertion follows from the component-by-component equation

$$x_i = r\left(\alpha_i - \frac{1}{2n}\right)$$

with $\alpha_i \in (0, 1/n)$ and the identity

$$\begin{aligned} \mathbf{x} &= \sum_{i=1}^n x_i \mathbf{e}_i \\ &= r \sum_{i=1}^n \alpha_i \mathbf{e}_i - \frac{r}{2n} \mathbf{1} \\ &= r \sum_{i=1}^n \alpha_i \left(\mathbf{e}_i - \frac{1}{2n} \mathbf{1} \right) - \left(1 - \sum_{i=1}^n \alpha_i \right) \frac{r}{2n} \mathbf{1}. \end{aligned}$$

The boundedness of $f(\mathbf{x})$ on J now follows from the inequalities

$$f\left(\sum_{i=0}^n \alpha_i \mathbf{v}_i\right) \leq \sum_{i=0}^n \alpha_i f(\mathbf{v}_i) \leq \max\{f(\mathbf{v}_0), \dots, f(\mathbf{v}_n)\}.$$

Without affecting its convexity, we now rescale and translate $f(\mathbf{x})$ so that $f(\mathbf{0}) = 0$ and $f(\mathbf{x}) \leq 1$ on J . We also rescale \mathbf{x} so that J contains the open ball $B(\mathbf{0}, 2)$. For any \mathbf{x} in $B(\mathbf{0}, 2)$, we have

$$0 = f(\mathbf{0}) \leq \frac{1}{2}f(\mathbf{x}) + \frac{1}{2}f(-\mathbf{x}).$$

It follows that $f(\mathbf{x})$ is bounded below by -1 on $B(\mathbf{0}, 2)$. The final step of the proof proceeds by choosing two distinct points \mathbf{x} and \mathbf{z} from the unit ball $B(\mathbf{0}, 1)$. If we define $\mathbf{w} = \mathbf{z} + t^{-1}(\mathbf{z} - \mathbf{x})$ with $t = \|\mathbf{z} - \mathbf{x}\|$, then $\mathbf{w} \in B(\mathbf{0}, 2)$,

$$\mathbf{z} = \frac{t}{1+t} \mathbf{w} + \frac{1}{1+t} \mathbf{x},$$

and

$$\begin{aligned} f(\mathbf{z}) - f(\mathbf{x}) &\leq \frac{t}{1+t} f(\mathbf{w}) + \frac{1}{1+t} f(\mathbf{x}) - f(\mathbf{x}) \\ &= \frac{t}{1+t} f(\mathbf{w}) - \frac{t}{1+t} f(\mathbf{x}) \\ &\leq \frac{2t}{1+t} \\ &\leq 2\|\mathbf{z} - \mathbf{x}\|. \end{aligned}$$

Switching the roles of \mathbf{x} and \mathbf{z} gives $|f(\mathbf{z}) - f(\mathbf{x})| \leq 2\|\mathbf{z} - \mathbf{x}\|$. This Lipschitz inequality establishes the continuity of $f(\mathbf{x})$ throughout $B(\mathbf{0}, 1)$. \blacksquare

We next turn to derivatives. The simplest place to start is forward directional derivatives. In the case of a convex function $f(x)$ defined on an interval (a, b) , we will prove the existence of one-sided derivatives by establishing the inequalities

$$\frac{f(y) - f(x)}{y - x} \leq \frac{f(z) - f(x)}{z - x} \leq \frac{f(z) - f(y)}{z - y} \quad (6.7)$$

for all points $x < y < z$ drawn from (a, b) . If we write

$$y = \frac{z-y}{z-x}x + \frac{y-x}{z-x}z,$$

then both of these inequalities are rearrangements of the inequality

$$f(y) \leq \frac{z-y}{z-x}f(x) + \frac{y-x}{z-x}f(z).$$

Careful examination of the inequalities (6.7) with relabeling of points as necessary leads to the conclusion that the slope

$$\frac{f(y) - f(x)}{y - x}$$

is bounded below and increasing in y for x fixed. Similarly, this same slope is bounded above and increasing in x for y fixed. It follows that both one-sided derivatives exist at y and satisfy

$$f'_-(y) = \lim_{x \uparrow y} \frac{f(y) - f(x)}{y - x} \leq \lim_{z \downarrow y} \frac{f(z) - f(y)}{z - y} = f'_+(y).$$

In view of the monotonicity properties of the slope, any number d between these two limits satisfies the supporting hyperplane inequalities

$$\begin{aligned} f(x) &\geq f(y) + d(x - y) \\ f(z) &\geq f(y) + d(z - y). \end{aligned}$$

Such a number is termed a subgradient. The existence of subgradients is closely tied to the fact that $f(x)$ is locally Lipschitz.

Our reasoning for convex functions defined on the real lines proves the existence of forward directional derivatives on higher-dimensional domains. Indeed, for any point \mathbf{y} in the domain of $f(\mathbf{x})$, all one must do is focus on the function $g(t) = f(\mathbf{y} + t\mathbf{v})$ of the nonnegative scalar t . To define the difference quotient of $g(t)$ at 0, the point $\mathbf{y} + t\mathbf{v}$ must belong to the domain of $f(\mathbf{x})$ for all t sufficiently small. This is certainly possible when \mathbf{y} occurs on the interior of the domain, but it is also possible for boundary points \mathbf{y} and directions \mathbf{v} that point into the domain of $f(\mathbf{x})$. In either case, the convexity of $f(\mathbf{x})$ carries over to $g(t)$. Furthermore, the directional derivative

$$d_{\mathbf{v}}f(\mathbf{y}) = \lim_{t \downarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = \lim_{t \downarrow 0} \frac{g(t) - g(0)}{t}.$$

is well defined and satisfies $d_{\mathbf{v}}f(\mathbf{y}) < \infty$. The value $d_{\mathbf{v}}f(\mathbf{y}) = -\infty$ can occur for boundary points \mathbf{y} but can be ruled out for interior points. At an interior point $g(t)$ is defined for t small and negative.

Finally, let us consider the integrability of a convex function $f(x)$ defined on a closed interval $[a, b]$. On the interior of the interval, $f(x)$ is continuous. Continuity can fail at the endpoints, but typically we can restore it by replacing $f(a)$ by $\lim_{x \rightarrow a} f(x)$ and $f(b)$ by $\lim_{x \rightarrow b} f(x)$. If we assume these limits are finite, then $f(x)$ is continuous throughout the interval and hence integrable. More surprising is the fact that the fundamental theorem of calculus applies. The right-hand and left-hand derivatives $f'_+(x)$ and $f'_-(x)$ exist throughout (a, b) , and

$$f(b) - f(a) = \int_a^b f'_+(x) dx = \int_a^b f'_-(x) dx. \quad (6.8)$$

Because $f'_-(x) \leq f'_+(x) \leq f'_-(y) \leq f'_+(y)$ when $x < y$, the interiors of the intervals $[f'_-(x), f'_+(x)]$ are disjoint. Choosing a rational number from each nonempty interior shows that the set of points where $f'_-(x) \neq f'_+(x)$ is countable. The value of an integral is insensitive to the value of its integrand at a countable number of points, and the discussion following Proposition 3.4.1 demonstrates that the fundamental theorem of calculus holds as stated in equation (6.8).

6.5 Minimization of Convex Functions

Optimization theory is much simpler for convex functions than for ordinary functions. The continuity and differentiability properties of convex functions certainly support this contention. Even more relevant are the following theoretical results.

Proposition 6.5.1 *Suppose that $f(\mathbf{y})$ is a convex function on the convex set $S \subset \mathbb{R}^n$. If \mathbf{x} is a local minimum of $f(\mathbf{y})$, then it is a global minimum of $f(\mathbf{y})$, and the set $\{\mathbf{y} \in S : f(\mathbf{y}) = f(\mathbf{x})\}$ is convex. If $f(\mathbf{y})$ is strictly convex and \mathbf{x} is a global minimum, then the solution set $\{\mathbf{y} \in S : f(\mathbf{y}) = f(\mathbf{x})\}$ consists of \mathbf{x} alone.*

Proof: If $f(\mathbf{y}) \leq f(\mathbf{x})$ and $f(\mathbf{z}) \leq f(\mathbf{x})$, then

$$f[\alpha \mathbf{y} + (1 - \alpha)\mathbf{z}] \leq \alpha f(\mathbf{y}) + (1 - \alpha)f(\mathbf{z}) \leq f(\mathbf{x}) \quad (6.9)$$

for any $\alpha \in [0, 1]$. This shows that the set $\{\mathbf{y} \in S : f(\mathbf{y}) \leq f(\mathbf{x})\}$ is convex. Now suppose that $f(\mathbf{y}) < f(\mathbf{x})$. Strict inequality then prevails between the extreme members of inequality (6.9) provided $\alpha > 0$. Taking $\mathbf{z} = \mathbf{x}$ and α close to 0 shows that \mathbf{x} cannot serve as a local minimum. This contradiction demonstrates that \mathbf{x} must be a global minimum. Finally, if $f(\mathbf{y})$ is strictly convex, then strict inequality holds in the first half of equality (6.9) for all $\alpha \in (0, 1)$. This leads to another contradiction when $\mathbf{y} = \mathbf{x} \neq \mathbf{z}$, and both are minimum points. ■

Example 6.5.1 *Piecewise Linear Functions*

The function $f(x) = |x|$ on the real line is piecewise linear. It attains its minimum of 0 at the point $x = 0$. The convex function $f(x) = \max\{1, |x|\}$ is also piecewise linear, but it attains its minimum throughout the interval $[-1, 1]$. In both cases the set $\{y : f(y) = \min_x f(x)\}$ is convex. In higher dimensions, the convex function $f(\mathbf{x}) = \max\{1, \|\mathbf{x}\|\}$ attains its minimum of 1 throughout the closed ball $\|\mathbf{x}\| \leq 1$. ■

Proposition 6.5.2 *Let $f(\mathbf{y})$ be a convex function defined on a convex set $S \subset \mathbb{R}^n$. A point $\mathbf{x} \in S$ furnishes a global minimum of $f(\mathbf{y})$ if and only if the forward directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ exists and is nonnegative for all tangent vectors $\mathbf{v} = \mathbf{z} - \mathbf{x}$ defined by $\mathbf{z} \in S$. In particular, a stationary point of $f(\mathbf{y})$ represents a global minimum.*

Proof: Suppose the condition holds and $\mathbf{z} \in S$. Then taking $t = 1$ and $\mathbf{v} = \mathbf{z} - \mathbf{x}$ in the inequality

$$\frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} \geq d_{\mathbf{v}}f(\mathbf{x})$$

shows that $f(\mathbf{z}) \geq f(\mathbf{x})$. Conversely, if \mathbf{x} represents the minimum, then the displayed difference quotient is nonnegative. Sending t to 0 now gives $d_{\mathbf{v}}f(\mathbf{x}) \geq 0$. ■

Example 6.5.2 *Minimum of y on $[0, \infty)$*

The convex function $f(y) = y$ has derivative $f'(y) = 1$. On the convex set $[0, \infty)$, we have $f'(0)(z - 0) = z \geq 0$ for any $z \in [0, \infty)$. Hence, 0 provides the minimum of y . Of course, this is consistent with the Lagrange multiplier rule $f'(0) - 1 = 0$. ■

Example 6.5.3 *The Obtuse Angle Criterion*

As a continuation of Propositions 6.2.2 and 6.2.3, define $f(\mathbf{y}) = \frac{1}{2}\|\mathbf{y} - \mathbf{x}\|^2$ for $\mathbf{x} \notin S$. If \mathbf{z} is the projection of \mathbf{x} onto S , then the necessary and sufficient condition for a minimum given by Proposition 6.5.2 reads

$$d_{\mathbf{v}}f(\mathbf{z}) = (\mathbf{z} - \mathbf{x})^* \mathbf{v} \geq 0$$

for every direction $\mathbf{v} = \mathbf{y} - \mathbf{z}$ defined by another point $\mathbf{y} \in S$. This inequality can be rephrased as the obtuse angle criterion $(\mathbf{x} - \mathbf{z})^*(\mathbf{y} - \mathbf{z}) \leq 0$ for all such \mathbf{y} . A simple diagram makes this conclusion visually obvious. ■

In Chap. 5 we found that the multiplier rule (5.1) is a necessary condition for a feasible point \mathbf{x} to be a local minimum of the objective function $f(\mathbf{y})$ subject to the constraints

$$\begin{aligned} g_i(\mathbf{y}) &= 0, & 1 \leq i \leq p \\ h_j(\mathbf{y}) &\leq 0, & 1 \leq j \leq q. \end{aligned}$$

In the presence of convexity, the multiplier rule is also a sufficient condition. We now revisit the entire issue under a combination of weaker and stronger hypotheses. The stronger hypotheses amount to: (a) $f(\mathbf{y})$ is a convex function, (b) the $g_i(\mathbf{y})$ are affine functions, and (c) the feasible region S is a convex set. Instead of assuming that $f(\mathbf{y})$ and the $h_j(\mathbf{y})$ are continuously differentiable, we now require them to be simply differentiable at the point of interest \mathbf{x} . Note that we do not require the $h_j(\mathbf{y})$ to be convex. Of course if they are, then S is automatically convex.

Proposition 6.5.3 *Under the above conditions, suppose the feasible point \mathbf{x} satisfies the multiplier rule (5.1) with $\lambda_0 = 1$. Then \mathbf{x} furnishes a global minimum of $f(\mathbf{y})$. Conversely, if \mathbf{x} is a minimum point satisfying the Mangasarian-Fromovitz constraint qualification, then the multiplier rule holds with $\lambda_0 = 1$.*

Proof: Suppose \mathbf{x} satisfies the multiplier rule. Take the inner product of the multiplier formula

$$\nabla f(\mathbf{x}) + \sum_{i=1}^p \lambda_i \nabla g_i(\mathbf{x}) + \sum_{j=1}^q \mu_j \nabla h_j(\mathbf{x}) = \mathbf{0}$$

with a vector $\mathbf{v} = \mathbf{z} - \mathbf{x}$ defined by a second feasible point \mathbf{z} . Because the equality constraint $g_i(\mathbf{y})$ is affine, $dg_i(\mathbf{x})\mathbf{v} = 0$. It follows that

$$df(\mathbf{x})\mathbf{v} = - \sum_{j=1}^q \mu_j dh_j(\mathbf{x})\mathbf{v}.$$

This representation puts us into position to apply Proposition 6.5.2. If some $h_j(\mathbf{x}) < 0$, then the multiplier $\mu_j = 0$. If $h_j(\mathbf{x}) = 0$, then the difference quotient inequality

$$\frac{h_j(\mathbf{x} + t\mathbf{v}) - h_j(\mathbf{x})}{t} \leq 0$$

holds for all $t \in (0, 1)$. Note here that the convexity of S subtly enters the argument. In any case, sending t to 0 produces $dh_j(\mathbf{x})\mathbf{v} \leq 0$. Because $\mu_j \geq 0$, we conclude that $df(\mathbf{x})\mathbf{v} \geq 0$, and this suffices to establish the claim that \mathbf{x} is a minimum point.

Proposition 14.7.1 proves the converse under relaxed differentiability assumptions. Here we limit ourselves to the case of no equality constraints ($p = 0$). In this setting we consider the convex function

$$m(\mathbf{y}) = \max\{f(\mathbf{y}) - f(\mathbf{x}), h_j(\mathbf{y}), 1 \leq j \leq q\}.$$

It is clear that the minimum of $m(\mathbf{y})$ occurs at \mathbf{x} because $m(\mathbf{x}) = 0$ and for all remaining \mathbf{y} either a constraint is violated or $f(\mathbf{y}) \geq f(\mathbf{x})$. Proposition 6.5.2 therefore implies $d_{\mathbf{v}}m(\mathbf{x}) \geq 0$ for all directions \mathbf{v} . Now let J

denote the index set $\{j : h_j(\mathbf{x}) = 0\}$. Since all of the functions defining $m(\mathbf{x})$ except the inactive inequality constraints achieve the maximum of 0, Example 4.4.4 yields the forward directional derivative

$$d_{\mathbf{v}}m(\mathbf{x}) = \max\{df(\mathbf{x})\mathbf{v}, dh_j(\mathbf{x})\mathbf{v}, 1 \leq i \leq p, j \in J\}.$$

Because $d_{\mathbf{v}}m(\mathbf{x}) \geq 0$ for all \mathbf{v} , Proposition 5.3.2 shows that there exist a convex combination of the relevant gradients with

$$\lambda_0 \nabla f(\mathbf{x}) + \sum_{j=1}^q \mu_j \nabla h_j(\mathbf{x}) = \mathbf{0}.$$

To eliminate the possibility $\lambda_0 = 0$, now invoke the argument in the last paragraph of the proof of Proposition 5.2.1. ■

Example 6.5.4 Slater's Constraint Qualification

The Mangasarian-Fromovitz constraint qualification is implied by a simpler condition called the Slater constraint qualification under affine equality constraints and convex inequality constraints. Slater's condition postulates the existence of a feasible point \mathbf{z} such that $h_j(\mathbf{z}) < 0$ for all j . If \mathbf{x} is a candidate minimum point and the row vectors $dg_i(\mathbf{x})$ are linearly independent, then the Mangasarian-Fromovitz constraint qualification involves finding a vector \mathbf{v} with $dg_i(\mathbf{x})\mathbf{v} = 0$ for all i and $dh_j(\mathbf{x})\mathbf{v} < 0$ for all inequality constraints active at \mathbf{x} . If $h_j(\mathbf{y})$ is active at \mathbf{x} , then the inequalities

$$0 > h_j(\mathbf{z}) - h_j(\mathbf{x}) \geq dh_j(\mathbf{x})(\mathbf{z} - \mathbf{x})$$

demonstrate that the vector $\mathbf{v} = \mathbf{z} - \mathbf{x}$ satisfies the Mangasarian-Fromovitz constraint qualification. Nothing in this argument depends on the objective function $f(\mathbf{y})$ being convex. ■

Example 6.5.5 Concave Constraints

Consider once again the general nonlinear programming problem covered by Proposition 5.2.1 with the proviso that the equality constraints $g_i(\mathbf{y})$ are affine and the inequality constraints $h_j(\mathbf{y})$ are concave rather than convex. One can eliminate the equality constraint $g_i(\mathbf{y}) = 0$ and replace it by the two inequality constraints $g_i(\mathbf{y}) \leq 0$ and $-g_i(\mathbf{y}) \leq 0$. Under this substitution all inequality constraints remain concave. In the multiplier rule (5.1) at a local minimum \mathbf{x} , concavity allows us to rule out the possibility that the multiplier λ_0 of $\nabla f(\mathbf{x})$ is 0. To prove this fact, assume that the first r inequality constraints are active at \mathbf{x} and the subsequent $q - r$ inequality constraints are inactive. Fortunately, Farkas' lemma poses a relevant dichotomy. Either

$$-\nabla f(\mathbf{x}) = \sum_{j=1}^r \mu_j \nabla h_j(\mathbf{x})$$

for nonnegative multipliers μ_j , or there is a vector \mathbf{w} with $-\mathbf{w}^* \nabla f(\mathbf{x}) > 0$ and $\mathbf{w}^* \nabla h_j(\mathbf{x}) \leq 0$ for all $j \leq r$. Suppose such a vector \mathbf{w} exists. If \mathbf{x} is an interior point of the common domain of $f(\mathbf{x})$ and the constraints $h_j(\mathbf{x})$, then on the one hand the inequality

$$h_j(\mathbf{x} + t\mathbf{w}) \leq h_j(\mathbf{x}) + tdh_j(\mathbf{x})\mathbf{w} \leq h_j(\mathbf{x}) \leq 0$$

shows that the point $\mathbf{x} + t\mathbf{w}$ is feasible for small $t > 0$. On the other hand,

$$f(\mathbf{x} + t\mathbf{w}) = f(\mathbf{x}) + tdf(\mathbf{x})\mathbf{w} + o(t) < f(\mathbf{x})$$

for small $t > 0$. This contradicts the assumption that \mathbf{x} is a local minimum, and the multiplier rule holds with $\lambda_0 = 1$. Linear programming is the most important application. Here the multiplier rule is a necessary and sufficient condition for a global minimum regardless of whether the active affine constraints are linearly independent. ■

Example 6.5.6 *Minimum of a Positive Definite Quadratic Function*

The quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ has gradient

$$\nabla f(\mathbf{x}) = \mathbf{A} \mathbf{x} + \mathbf{b}$$

for \mathbf{A} symmetric. Assuming that \mathbf{A} is positive definite and affine equality constraints are present, Proposition 6.5.3 demonstrates that the candidate minimum point identified in Example 5.2.6 furnishes the global minimum of $f(\mathbf{x})$. When \mathbf{A} is merely positive semidefinite, Example 6.5.5 shows that the multiplier rule with $\lambda_0 = 1$ is still a necessary and sufficient condition for a minimum. ■

Example 6.5.7 *Maximum Likelihood for the Multivariate Normal*

The sample mean and sample variance

$$\begin{aligned} \bar{\mathbf{y}} &= \frac{1}{k} \sum_{j=1}^k \mathbf{y}_j \\ \mathbf{S} &= \frac{1}{k} \sum_{j=1}^k (\mathbf{y}_j - \bar{\mathbf{y}})(\mathbf{y}_j - \bar{\mathbf{y}})^* \end{aligned}$$

are also the maximum likelihood estimates of the theoretical mean $\boldsymbol{\mu}$ and theoretical variance $\boldsymbol{\Omega}$ of a random sample $\mathbf{y}_1, \dots, \mathbf{y}_k$ from a multivariate normal distribution. (See Appendix A.2 for a review of the multivariate normal.) To prove this fact, we first note that maximizing the loglikelihood function

$$-\frac{k}{2} \ln \det \boldsymbol{\Omega} - \frac{1}{2} \sum_{j=1}^k (\mathbf{y}_j - \boldsymbol{\mu})^* \boldsymbol{\Omega}^{-1} (\mathbf{y}_j - \boldsymbol{\mu})$$

$$\begin{aligned}
&= -\frac{k}{2} \ln \det \boldsymbol{\Omega} - \frac{k}{2} \boldsymbol{\mu}^* \boldsymbol{\Omega}^{-1} \boldsymbol{\mu} + \left(\sum_{j=1}^k \mathbf{y}_j \right)^* \boldsymbol{\Omega}^{-1} \boldsymbol{\mu} - \frac{1}{2} \sum_{j=1}^k \mathbf{y}_j^* \boldsymbol{\Omega}^{-1} \mathbf{y}_j \\
&= -\frac{k}{2} \ln \det \boldsymbol{\Omega} - \frac{1}{2} \operatorname{tr} \left[\boldsymbol{\Omega}^{-1} \sum_{j=1}^k (\mathbf{y}_j - \boldsymbol{\mu})(\mathbf{y}_j - \boldsymbol{\mu})^* \right]
\end{aligned}$$

with respect to $\boldsymbol{\mu}$ constitutes a special case of the previous example with $\mathbf{A} = k\boldsymbol{\Omega}^{-1}$ and $\mathbf{b} = -\boldsymbol{\Omega}^{-1} \sum_{j=1}^k \mathbf{y}_j$. This leads to the same estimate $\hat{\boldsymbol{\mu}} = \bar{\mathbf{y}}$ regardless of the value of $\boldsymbol{\Omega}$. Once we fix $\hat{\boldsymbol{\mu}}$, we left with the problem of estimating $\boldsymbol{\Omega}$. Fortunately, this is a special case of the problem treated at the end of Example 4.7.6. The point $\hat{\boldsymbol{\Omega}} = \mathbf{S}$ corresponds to a maximum because the function $-\ln \det \boldsymbol{\Omega}$ is log-concave in $\boldsymbol{\Omega}^{-1}$ and the function $\operatorname{tr}(\boldsymbol{\Omega}^{-1}\mathbf{S})$ is linear in $\boldsymbol{\Omega}^{-1}$. Here we implicitly assume that \mathbf{S} is invertible.

Alternatively, we can estimate $\boldsymbol{\Omega}$ by exploiting the Cholesky decompositions $\boldsymbol{\Omega} = \mathbf{L}\mathbf{L}^*$ and $\mathbf{S} = \mathbf{M}\mathbf{M}^*$. (See Problems 37 and 38 for a development of the Cholesky decomposition of a positive definite matrix.) In view of the identities $\boldsymbol{\Omega}^{-1} = (\mathbf{L}^{-1})^* \mathbf{L}^{-1}$ and $\det \boldsymbol{\Omega} = (\det \mathbf{L})^2$, the loglikelihood becomes

$$\begin{aligned}
&k \ln \det \mathbf{L}^{-1} - \frac{k}{2} \operatorname{tr} \left[(\mathbf{L}^{-1})^* \mathbf{L}^{-1} \mathbf{M}\mathbf{M}^* \right] \\
&= k \ln \det (\mathbf{L}^{-1}\mathbf{M}) - \frac{k}{2} \operatorname{tr} \left[(\mathbf{L}^{-1}\mathbf{M})(\mathbf{L}^{-1}\mathbf{M})^* \right] - k \ln \det \mathbf{M}
\end{aligned}$$

using the cyclic permutation property of the matrix trace function. Because products and inverses of lower triangular matrices are lower triangular, the matrix $\mathbf{R} = \mathbf{L}^{-1}\mathbf{M}$ ranges over the set of lower triangular matrices with positive diagonal entries as \mathbf{L} ranges over the same set. This permits us to reparameterize and estimate $\mathbf{R} = (r_{ij})$ instead of \mathbf{L} . Up to an irrelevant additive constant, the loglikelihood reduces to

$$k \ln \det \mathbf{R} - \frac{k}{2} \operatorname{tr}(\mathbf{R}\mathbf{R}^*) = k \sum_i \ln r_{ii} - \frac{k}{2} \sum_i \sum_{j=1}^i r_{ij}^2.$$

Clearly, this is maximized by taking $r_{ij} = 0$ for $j \neq i$. Differentiation of the concave function $k \ln r_{ii} - \frac{k}{2} r_{ii}^2$ shows that it is maximized by taking $r_{ii} = 1$. In other words, the maximum likelihood estimator $\hat{\mathbf{R}}$ is the identity matrix \mathbf{I} . This implies that $\hat{\mathbf{L}} = \mathbf{M}$ and consequently that $\hat{\boldsymbol{\Omega}} = \mathbf{S}$. ■

Example 6.5.8 Geometric Programming

The function

$$f(\mathbf{t}) = \sum_{i=1}^j c_i \prod_{k=1}^n t_k^{\beta_{ik}}$$

is called a posynomial if all components t_1, \dots, t_n of the argument \mathbf{t} and all coefficients c_1, \dots, c_j are positive. The powers β_{ik} may be positive, negative, or zero. For instance, $t_1^{-1} + 2t_1^3t_2^{-2}$ is a posynomial on \mathbb{R}^2 . Geometric programming deals with the minimization of a posynomial $f(\mathbf{t})$ subject to posynomial inequality constraints of the form $h_j(\mathbf{t}) \leq 1$ for $1 \leq j \leq q$.

Better understanding of geometric programming can be achieved by making the change of variables $t_k = e^{x_k}$. This eliminates the constraint $t_k > 0$ and shows that

$$g(\mathbf{x}) = \sum_{i=1}^j c_i \prod_{k=1}^n t_k^{\beta_{ik}} = \sum_{i=1}^j c_i e^{\beta_i^* \mathbf{x}}$$

is log-convex in the transformed parameters. The reparameterized constraint functions are likewise log-convex and define a convex feasible region S . If the vectors β_1, \dots, β_j span \mathbb{R}^n , then the expression

$$d^2g(\mathbf{x}) = \sum_{i=1}^j c_i e^{\beta_i^* \mathbf{x}} \beta_i \beta_i^*$$

for the second differential proves that $g(\mathbf{x})$ is strictly convex. It follows that if $g(\mathbf{x})$ possesses a minimum, then it is achieved at a single point. ■

Example 6.5.9 *Quasi-Convexity*

If $f(\mathbf{x})$ is convex and $g(z)$ is an increasing function of the real variable z , then the inequality

$$\begin{aligned} f[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] &\leq \alpha f(\mathbf{x}) + (1 - \alpha)f(\mathbf{y}) \\ &\leq \max\{f(\mathbf{x}), f(\mathbf{y})\} \end{aligned}$$

implies that the function $h(\mathbf{x}) = g \circ f(\mathbf{x})$ satisfies

$$h[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] \leq \max\{h(\mathbf{x}), h(\mathbf{y})\} \tag{6.10}$$

for any $\alpha \in [0, 1]$. Satisfaction of inequality (6.10) is sometimes taken as the definition of quasi-convexity for an arbitrary function $h(\mathbf{x})$. If $f(\mathbf{x})$ is strictly convex and $g(z)$ is strictly increasing, then strict inequality prevails in inequality (6.10) when $\alpha \in (0, 1)$ and $\mathbf{y} \neq \mathbf{x}$. Once again this implication can be turned into a definition. The importance of strict quasi-convexity lies in the fact that a strictly quasi-convex function possesses at most one local minimum, and if a local minimum exists, then it is necessarily the global minimum.

Similar considerations apply to concave and quasi-concave functions. For example, the function $h(x) = e^{-(x-\mu)^2}$ is strictly quasi-concave because it is the composition of the strictly increasing function $g(y) = e^y$ with the strictly concave function $f(x) = -(x-\mu)^2$. It is clear that $h(x)$ has a global maximum at $x = \mu$. ■

6.6 Moment Inequalities

In this section we assume that readers have a good grasp of probability theory. For those with limited background, most of the material can be comprehended by restricting attention to discrete random variables.

Inequalities give important information about the magnitude of probabilities and expectations without requiring their exact calculation. The Cauchy-Schwarz inequality $|E(XY)| \leq E(X^2)^{1/2} E(Y^2)^{1/2}$ is one of the most useful of the classical inequalities. (The reader should check that this is just a disguised form of the Cauchy-Schwarz inequality of Chap. 1 applied to random variables X and Y .) It is also one of the easiest to remember because it is equivalent to the fact that a correlation coefficient must lie on the interval $[-1, 1]$. Equality occurs in the Cauchy-Schwarz inequality if and only if X is proportional to Y or vice versa.

Markov's inequality is another widely applied bound. Let $g(x)$ be a non-negative, increasing function, and let X be a random variable such that $g(X)$ has finite expectation. Then Markov's inequality

$$\Pr(X \geq c) \leq \frac{E[g(X)]}{g(c)}$$

holds for any constant c for which $g(c) > 0$ and follows logically by taking expectations in the inequality $g(c)1_{\{X \geq c\}} \leq g(X)$. Chebyshev's inequality is the special case of Markov's inequality with $g(x) = x^2$ applied to the random variable $|X - E(X)|$. Chebyshev's inequality reads

$$\Pr[|X - E(X)| \geq c] \leq \frac{\text{Var}(X)}{c^2}.$$

In large deviation theory, we take $g(x) = e^{tx}$ and $c > 0$ and choose $t > 0$ to minimize the right-hand side of the inequality $\Pr(X \geq c) \leq e^{-ct} E(e^{tX})$ involving the moment generating function of X . As an example, suppose X follows a standard normal distribution. The moment generating function $e^{t^2/2}$ of X is derived by a minor variation of the argument given in Appendix A.1 for the characteristic function of X . The large deviation inequality

$$\Pr[X \geq c] \leq \inf_t e^{-ct} e^{t^2/2} = e^{-c^2/2}$$

is called a Chernoff bound. Problem 43 discusses another typical Chernoff bound.

Our next example involves a nontrivial application of Chebyshev's inequality. In preparation for the example, we recall that a binomially distributed random variable S_n has distribution

$$\Pr(S_n = k) = \binom{n}{k} x^k (1-x)^{n-k}.$$

Here S_n is interpreted as the number of successes in n independent trials with success probability x per trial [89]. The mean and variance of S_n are $E(S_n) = nx$ and $\text{Var}(S_n) = nx(1-x)$.

Example 6.6.1 *Weierstrass's Approximation Theorem*

Weierstrass showed that a continuous function $f(x)$ on $[0, 1]$ can be uniformly approximated to any desired degree of accuracy by a polynomial. Bernstein's lovely proof of this fact relies on applying Chebyshev's inequality to the random variable S_n/n derived from the binomial random variable S_n just discussed. The corresponding candidate polynomial is defined by the expectation

$$E \left[f \left(\frac{S_n}{n} \right) \right] = \sum_{k=0}^n f \left(\frac{k}{n} \right) \binom{n}{k} x^k (1-x)^{n-k}.$$

Note that $E(S_n/n) = x$ and

$$\text{Var} \left(\frac{S_n}{n} \right) = \frac{x(1-x)}{n} \leq \frac{1}{4n}.$$

Now given an arbitrary $\epsilon > 0$, one can find by the uniform continuity of $f(x)$ a $\delta > 0$ such that $|f(u) - f(v)| < \epsilon$ whenever $|u - v| < \delta$. If $\|f\|_\infty = \sup |f(x)|$ on $[0, 1]$, then Chebyshev's inequality implies

$$\begin{aligned} & \left| E \left[f \left(\frac{S_n}{n} \right) \right] - f(x) \right| \\ & \leq E \left[\left| f \left(\frac{S_n}{n} \right) - f(x) \right| \right] \\ & \leq \epsilon \Pr \left(\left| \frac{S_n}{n} - x \right| < \delta \right) + 2\|f\|_\infty \Pr \left(\left| \frac{S_n}{n} - x \right| \geq \delta \right) \\ & \leq \epsilon + \frac{2\|f\|_\infty x(1-x)}{n\delta^2} \\ & \leq \epsilon + \frac{\|f\|_\infty}{2n\delta^2}. \end{aligned}$$

Taking $n \geq \|f\|_\infty / (2\epsilon\delta^2)$ then gives $\left| E \left[f \left(\frac{S_n}{n} \right) \right] - f(x) \right| \leq 2\epsilon$ regardless of the chosen $x \in [0, 1]$. ■

Proposition 6.6.1 (Jensen) *Let the values of the random variable W be confined to the possibly infinite interval (a, b) . If $h(w)$ is convex on (a, b) , then $E[h(W)] \geq h[E(W)]$, provided both expectations exist. For a strictly convex function $h(w)$, equality holds in Jensen's inequality if and only if $W = E(W)$ almost surely.*

Proof: For the sake of simplicity, assume that $h(w)$ is differentiable at the point $v = E(W)$. Then Jensen's inequality follows from Proposition 6.3.1

after taking expectations in the inequality

$$h(W) \geq h(v) + dh(v)(W - v). \quad (6.11)$$

If $h(w)$ is strictly convex, and W is not constant, then inequality (6.11) is strict with positive probability. Hence, strict inequality prevails in Jensen's inequality. As we will see later, the differentiability assumption on $h(w)$ can be relaxed by substituting a subgradient for the gradient. ■

Jensen's inequality is the key to a host of other inequalities. Here is one important example.

Example 6.6.2 *Schlömilch's Inequality for Weighted Means*

If X is a nonnegative random variable, then we define the weighted mean function $M(p) = E(X^p)^{\frac{1}{p}}$. For the sake of argument, we assume that $M(p)$ exists and is finite for all real p . Typical values of $M(p)$ are $M(1) = E(X)$ and $M(-1) = 1/E(X^{-1})$. To make $M(p)$ continuous at $p = 0$, it turns out that we should set $M(0) = e^{E(\ln X)}$. The reader is asked to check this fact in Problem 50. Here we are more concerned with proving Schlömilch's assertion that $M(p)$ is an increasing function of p . If $0 < p < q$, then the function $x \mapsto x^{q/p}$ is convex, and Jensen's inequality says

$$E(X^p)^{q/p} \leq E(X^q).$$

Taking the q th root of both sides of this inequality yields $M(p) \leq M(q)$. On the other hand if $p < q < 0$, then the function $x \mapsto x^{q/p}$ is concave, and Jensen's inequality says

$$E(X^p)^{q/p} \geq E(X^q).$$

Taking the q th root reverses the inequality and again yields $M(p) \leq M(q)$. When either p or q is 0, we have to change tactics. One approach is to invoke the continuity of $M(p)$ at $p = 0$. Another approach is to exploit the concavity of $\ln x$. Jensen's inequality now gives

$$E(\ln X^p) \leq \ln E(X^p),$$

which on exponentiation becomes

$$e^{pE(\ln X)} \leq E(X^p).$$

If $p > 0$, then taking the p th root produces

$$M(0) = e^{E(\ln X)} \leq E(X^p)^{\frac{1}{p}},$$

and if $p < 0$, then taking the p th root produces the opposite inequality

$$M(0) = e^{E(\ln X)} \geq E(X^p)^{\frac{1}{p}}.$$

When the random variable X is confined to the space $\{1, \dots, n\}$ equipped with the uniform probabilities $p_i = 1/n$, Schlömilch's inequalities for the values $p = -1, 0$, and 1 reduce to the classical inequalities

$$\frac{1}{\frac{1}{n}\left(\frac{1}{x_1} + \dots + \frac{1}{x_n}\right)} \leq \left(x_1 \cdots x_n\right)^{\frac{1}{n}} \leq \frac{1}{n}\left(x_1 + \dots + x_n\right)$$

relating the harmonic, geometric, and arithmetic means. ■

Example 6.6.3 Hölder's Inequality

Consider two random variables X and Y and two numbers $p > 1$ and $q > 1$ such that $p^{-1} + q^{-1} = 1$. Then Hölder's inequality

$$|E(XY)| \leq E(|X|^p)^{\frac{1}{p}} E(|Y|^q)^{\frac{1}{q}} \quad (6.12)$$

generalizes the Cauchy-Schwarz inequality whenever the indicated expectations on its right exist. To prove (6.12), it clearly suffices to assume that X and Y are nonnegative. It also suffices to take $E(X^p) = E(Y^q) = 1$ once we divide the left-hand side of (6.12) by its right-hand side. To complete the proof, substitute the random variables X and Y for the scalars x and y in Young's inequality (1.3) and take expectations. ■

6.7 Problems

1. Suppose S and T are nonempty closed convex sets with empty intersection. Prove that there exists a unit vector \mathbf{v} such that

$$\sup_{\mathbf{x} \in S} \mathbf{v}^* \mathbf{x} \leq \inf_{\mathbf{y} \in T} \mathbf{v}^* \mathbf{y}.$$

If either S or T is bounded, then demonstrate further that strict inequality prevails in this inequality. (Hints: The set $S - T$ is convex and does not contain $\mathbf{0}$. If S is compact, then $S - T$ is closed.)

2. Let C be a convex set situated in \mathbb{R}^{m+n} . Show that the projected set $\{\mathbf{x} \in \mathbb{R}^m : (\mathbf{x}, \mathbf{y}) \in C \text{ for some } \mathbf{y} \in \mathbb{R}^n\}$ is convex.
3. Demonstrate that the set

$$S = \{\mathbf{x} \in \mathbb{R}^2 : x_1 x_2 \geq 1, x_1 \geq 0, x_2 \geq 0\}$$

is closed and convex. Further show that $P(S)$ is convex but not closed, where $P\mathbf{x} = x_1$ denotes projection onto the first coordinate of \mathbf{x} .

4. The function $P(\mathbf{x}, t) = t^{-1}\mathbf{x}$ from $\mathbb{R}^n \times (0, \infty)$ to \mathbb{R}^n is called the perspective map. Show that the image $P(C)$ and inverse image $P^{-1}(D)$ of convex sets C and D are convex. (Hint: Prove that $P(\mathbf{x}, t)$ maps line segments onto line segment.)

5. The topological notions of open, closed, and bounded can be recast for a convex set $C \subset \mathbb{R}^n$ [172]. Validate the necessary and sufficient conditions in each of the following assertions:
- C is open if and only if for every $\mathbf{x} \in C$ and $\mathbf{y} \in \mathbb{R}^n$, the point $\mathbf{x} + t\mathbf{y}$ belongs to C for all sufficiently small positive t .
 - C is closed if and only if whenever C contains the open segment $\{\alpha\mathbf{x} + (1 - \alpha)\mathbf{y} : \alpha \in (0, 1)\}$, it also contains the endpoints \mathbf{x} and \mathbf{y} of the segment.
 - C is bounded if and only if it contains no ray $\{\mathbf{x} + t\mathbf{y} : t \in [0, \infty)\}$.
6. Demonstrate that the convex hull of an open set is open.
7. The Gauss-Lucas theorem says that the roots of the derivative $p'(z)$ of a polynomial $p(z)$ are contained in the convex hull of the roots of $p(z)$. Prove this claim by exploiting the expansion

$$0 = \frac{p'(y)}{p(y)} = \sum_i \frac{1}{y - z_i} = \sum_i \frac{\bar{y} - \bar{z}_i}{\|y - z_i\|^2},$$

where y is a root of $p'(z)$ differing from each of the roots z_i of $p(z)$, and the overbar sign denotes complex conjugation.

8. Deduce Farkas' lemma from Proposition 5.3.2.
9. On which intervals are the following functions convex: e^x , e^{-x} , x^n for n an integer, $|x|^p$ for $p \geq 1$, $\sqrt{1 + x^2}$, $x \ln x$, and $\cosh x$? On these intervals, which functions are log-convex?
10. Demonstrate that the function $f(x) = x^n - na \ln x$ is convex on $(0, \infty)$ for any positive real number a and nonnegative integer n . Where does its minimum occur?
11. Show that Riemann's zeta function

$$\zeta(s) = \sum_{n=1}^{\infty} \frac{1}{n^s}$$

is log-convex for $s > 1$.

12. Show that the function

$$f(x) = \begin{cases} \frac{1 - e^{-xt}}{x} & x \neq 0 \\ t & x = 0 \end{cases}$$

is log-convex for $t > 0$ fixed. (Hints: Use either the second derivative test, or express $f(x)$ as the integral

$$f(x) = \int_0^t e^{-xs} ds,$$

and use the closure properties of log-convex functions.)

13. Use Proposition 6.3.1 to prove the Cauchy-Schwarz inequality.
14. Prove the strict convexity assertions of Proposition 6.3.1.
15. Prove parts (b), (c), (d), and (e) of Proposition 6.3.3.
16. Prove the unproved assertions of Proposition 6.3.4.
17. Let $f(x)$ and $g(x)$ be positive functions defined on an interval of the real line. Prove that:
 - (a) If $f(x)$ and $g(x)$ are both convex and increasing or both convex and decreasing, then the product $f(x)g(x)$ is convex.
 - (b) If $f(x)$ and $g(x)$ are both concave, one is increasing, and the other is decreasing, then the product $f(x)g(x)$ is concave.
 - (c) If $f(x)$ is convex and increasing and $g(x)$ is concave and decreasing, then the ratio $f(x)/g(x)$ is convex.

Here “increasing” means nondecreasing and similarly for “decreasing.” Your proofs should not assume that either $f(x)$ or $g(x)$ is differentiable.

18. Let $f(x)$ be a convex function on the real line that is bounded above. Demonstrate that $f(x)$ is constant.
19. Suppose $f(\mathbf{x}, \mathbf{y})$ is jointly convex in its two arguments and C is a convex set. Show that the function

$$g(\mathbf{x}) = \inf_{\mathbf{y} \in C} f(\mathbf{x}, \mathbf{y})$$

is convex. Assume here that the infimum is finite. As a special case, demonstrate that the distance function $\text{dist}(\mathbf{x}, C) = \inf_{\mathbf{y} \in C} \|\mathbf{y} - \mathbf{x}\|_*$ is convex for any norm $\|\cdot\|_*$. (Hint: For \mathbf{x}_1 and \mathbf{x}_2 and $\epsilon > 0$ choose points \mathbf{y}_1 and \mathbf{y}_2 in C satisfying $f(\mathbf{x}_i, \mathbf{y}_i) \leq g(\mathbf{x}_i) + \epsilon$.)

20. Let the function $f(x, y)$ be convex in y for each fixed x . Replace x by a random variable X and take expectations. Show that $E[f(X, y)]$ is convex in y , assuming the expectations in question exist. For n an even positive integer with $E(X^n) < \infty$, use this result to prove that $y \mapsto E[(X - y)^n]$ is convex in y .
21. Suppose $f(x)$ is absolutely integrable over every compact interval $[a, b]$ and $\phi(x)$ is nonnegative, bounded, vanishes outside a symmetric interval $[-c, c]$, and has total mass $\int \phi(x) dx = 1$. The convolution of f and $\phi(x)$ is defined by

$$f * \phi(x) = \int f(x - y)\phi(y) dy = \int f(y)\phi(x - y) dy.$$

Prove the following claims:

- (a) $f * \phi(x)$ is continuous if $\phi(x)$ is continuous.
- (b) $f * \phi(x)$ is k times continuously differentiable if $\phi(x)$ is k times continuously differentiable.
- (c) $f * \phi(x)$ is nonnegative if $f(x)$ is nonnegative.
- (d) $f * \phi(x)$ is increasing (decreasing) if $f(x)$ is increasing (decreasing).
- (e) $f * \phi(x)$ is convex (concave) if $f(x)$ is convex (concave).
- (f) $f * \phi(x)$ is log-convex if $f(x)$ is log-convex.

Define $\phi_n(s) = n\phi(nx)$ and prove as well that $f * \phi_n(x)$ converges to $f(x)$ at every point of continuity of $f(x)$. This convergence is uniform on every compact interval on which $f(x)$ is continuous. These facts imply that convexity properties established by differentiation carry over to convex functions in general. Can you supply any examples?

- 22. Suppose the polynomial $p(x)$ has only real roots. Show that $1/p(x)$ is log-convex on any interval where $p(x)$ is positive.
- 23. Demonstrate that the Kullback-Leibler (cross-entropy) distance

$$f(\mathbf{x}) = x_1 \ln \frac{x_1}{x_2} + x_2 - x_1$$

is convex on the set $\{\mathbf{x} = (x_1, x_2) : x_1 > 0, x_2 > 0\}$.

- 24. Show that the function $f(\mathbf{x}) = x_1^2 + x_2^4$ on \mathbb{R}^2 is strictly convex even though $d^2 f(\mathbf{x})$ is singular along the line $x_2 = 0$.
- 25. Prove that:

- (a) $f(\mathbf{x}) = x_1^2/x_2$ is convex for x_2 positive,
- (b) $f(\mathbf{x}) = \left(\prod_{i=1}^n x_i\right)^{1/n}$ is concave when all $x_i > 0$,
- (c) $f(\mathbf{x}) = \sum_{i=1}^m x_{[i]}$ is concave, where $x_{[1]} \leq \dots \leq x_{[m]}$ are the order statistics of $\mathbf{x} = (x_1, \dots, x_n)^*$,
- (d) $f(\mathbf{x}) = \left(\sum_{i=1}^n \sqrt{x_i}\right)^{-1}$ is convex when all $x_i \geq 0$ and at least one is positive.

(Hints: Use the second derivative test for (a), (b), and (d). Write

$$f(\mathbf{x}) = \min_{i_1 < \dots < i_m} [x_{i_1} + \dots + x_{i_m}]$$

for (c).)

26. Let $f(x)$ be a continuous function on the real line satisfying

$$f\left[\frac{1}{2}(x+y)\right] \leq \frac{1}{2}f(x) + \frac{1}{2}f(y).$$

Demonstrate that $f(x)$ is convex.

27. If $f(x)$ is a nondecreasing function on the interval $[a, b]$, then show that $g(x) = \int_a^x f(y)dy$ is a convex function on $[a, b]$.
28. The Bohr-Mollerup theorem asserts that $\Gamma(z)$ is the only log-convex function on the interval $(0, \infty)$ that satisfies $\Gamma(1) = 1$ and the factorial identity $\Gamma(z+1) = z\Gamma(z)$ for all z . We have seen that $\Gamma(z)$ has these properties. Prove conversely that any function $G(z)$ with these properties coincides with $\Gamma(z)$. (Hints: Check the inequalities

$$\begin{aligned} G(n+z) &\leq G(n)^{1-z}G(n)^zn^z = (n-1)!n^z \\ G(n+1) &\leq G(n+z)^zG(n+1+z)^{1-z} = G(n+z)(n+z)^{1-z} \end{aligned}$$

for all positive integers n and real numbers $z \in (0, 1)$. These in turn yield the inequalities

$$\frac{n!n^z}{z(z+1)\cdots(z+n)} \leq G(z) \leq \frac{n!n^z}{z(z+1)\cdots(z+n)} \cdot \frac{z+n}{n}.$$

Taking limits on n shows that $G(z)$ equals Gauss's infinite product expansion of $\Gamma(z)$. Note that this proof simultaneously validates Gauss's expansion (6.6).

29. Let $f(\mathbf{x})$ be a real-valued differentiable function on \mathbb{R}^n . If $f(\mathbf{x})$ is strictly convex, prove that $df(\mathbf{x}) = df(\mathbf{y})$ if and only if $\mathbf{x} = \mathbf{y}$.
30. Suppose $f(\mathbf{x})$ is convex on \mathbb{R}^n and $f(\mathbf{y}) = 0$. Prove that $f(\mathbf{x})^2$ is differentiable at \mathbf{y} with differential $\mathbf{0}^*$. (Hint: Invoke Proposition 6.4.1 and Fréchet's definition of differentiability.)
31. Let $f(\mathbf{x})$ be a convex differentiable function on \mathbb{R}^n . Show that the function

$$g(\mathbf{x}) = f(\mathbf{x}) + \epsilon\|\mathbf{x}\|^2$$

is strictly convex for $\epsilon > 0$ and that the set $\{\mathbf{x} \in \mathbb{R}^n : g(\mathbf{x}) \leq g(\mathbf{y})\}$ is compact for any \mathbf{y} .

32. A square matrix \mathbf{M} is said to be primitive if its entries are nonnegative and all of the entries of some power \mathbf{M}^p of \mathbf{M} are positive. The Perron-Frobenius theorem [136, 150, 234] asserts that a primitive matrix \mathbf{M} has a dominant eigenvalue $\lambda > 0$ possessing unique right and

left eigenvectors \mathbf{u} and \mathbf{v}^* with positive entries. Furthermore, if we choose \mathbf{u} and \mathbf{v}^* so that $\mathbf{v}^*\mathbf{u} = 1$, then

$$\lim_{m \rightarrow \infty} \lambda^{-m} \mathbf{M}^m = \mathbf{u}\mathbf{v}^*.$$

From this limit deduce that λ is log-convex in \mathbf{x} if the entries of \mathbf{M} are either 0 or log-convex functions of \mathbf{x} .

33. Consider minimizing the quadratic function $f(\mathbf{y}) = \frac{1}{2}\mathbf{y}^*\mathbf{A}\mathbf{y} + \mathbf{b}^*\mathbf{y} + c$ subject to the vector constraint $\mathbf{y} \geq \mathbf{0}$ for a positive semidefinite matrix \mathbf{A} . Show that the three conditions $\mathbf{A}\mathbf{x} + \mathbf{b} \geq \mathbf{0}$, $\mathbf{x} \geq \mathbf{0}$, and $\mathbf{x}^*\mathbf{A}\mathbf{x} + \mathbf{b}^*\mathbf{x} = 0$ are necessary and sufficient for \mathbf{x} to represent a minimum.
34. Let C be a convex set. Proposition 6.5.2 declares that a point $\mathbf{x} \in C$ minimizes a convex function $f(\mathbf{y})$ on C provided $d_{\mathbf{v}}f(\mathbf{x}) \geq 0$ for every tangent vector $\mathbf{v} = \mathbf{y} - \mathbf{x}$ constructed from a point $\mathbf{y} \in C$. If C is a convex cone, and $f(\mathbf{y})$ is convex and differentiable, then show that this condition is equivalent to the conditions $df(\mathbf{x})\mathbf{x} = 0$ and $df(\mathbf{x})\mathbf{y} \geq 0$ for every $\mathbf{y} \in C$.
35. The posynomial $f(\mathbf{x}) = \prod_{i=1}^n x_i^{\alpha_i}$ achieves a unique maximum on the unit simplex $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} \geq \mathbf{0}, \sum_{i=1}^n x_i = 1\}$ whenever the powers α_i are positive. Find this maximum, and show that it is global. (Hint: Minimize $-\ln f(\mathbf{x})$.)
36. Show that the functions $\sqrt{|x|}$, $\ln x$, and $\lfloor x \rfloor$ are quasi-convex. To the extent possible, state and prove the quasi-convex analogues of the convex closure properties covered in Proposition 6.3.3.
37. Let \mathbf{A} be an $n \times n$ positive definite matrix. The Cholesky decomposition \mathbf{B} of \mathbf{A} is a lower-triangular matrix with positive diagonal entries such that $\mathbf{A} = \mathbf{B}\mathbf{B}^*$. To prove that such a decomposition exists we can argue by induction. Why is the case of a 1×1 matrix trivial? Now suppose \mathbf{A} is partitioned as

$$\mathbf{A} = \begin{pmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{pmatrix}.$$

Applying the induction hypothesis, there exist matrices \mathbf{C}_{11} and \mathbf{D}_{22} such that

$$\begin{aligned} \mathbf{C}_{11}\mathbf{C}_{11}^* &= \mathbf{A}_{11} \\ \mathbf{D}_{22}\mathbf{D}_{22}^* &= \mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}. \end{aligned}$$

Prove that

$$\mathbf{B} = \begin{pmatrix} \mathbf{C}_{11} & \mathbf{0} \\ \mathbf{A}_{21}(\mathbf{C}_{11}^*)^{-1} & \mathbf{D}_{22} \end{pmatrix}$$

gives the desired decomposition. Extend this argument to show that \mathbf{B} is uniquely determined.

38. Continuing Problem 37, show that one can compute the Cholesky decomposition $\mathbf{B} = (b_{ij})$ of $\mathbf{A} = (a_{ij})$ by the recurrence relations

$$b_{jj} = \sqrt{a_{jj} - \sum_{k=1}^{j-1} b_{jk}^2}$$

$$b_{ij} = \frac{a_{ij} - \sum_{k=1}^{j-1} b_{ik}b_{jk}}{b_{jj}}, \quad i > j$$

for columns $j = 1, j = 2$, and so forth until column $j = n$. How can you compute $\det \mathbf{A}$ in terms of the entries of \mathbf{B} ?

39. Prove that the set of lower triangular matrices with positive diagonal entries is closed under matrix multiplication and matrix inversion.
40. Let X_1, \dots, X_n be n independent random variables from a common distributional family. Suppose the variance $\sigma^2(\mu)$ of a generic member of this family is a function of the mean μ . Now consider the sum $S = X_1 + \dots + X_n$. If the mean $\omega = E(S)$ is fixed, it is of some interest to determine whether taking $E(X_i) = \mu_i = \omega/n$ minimizes or maximizes $\text{Var}(S)$. Show that the minimum occurs when $\sigma^2(\mu)$ is convex in μ and the maximum occurs when $\sigma^2(\mu)$ is concave in μ . What do you deduce in the special cases where the family is binomial, Poisson, and exponential [201]?
41. If the random variable X has values in the interval $[a, b]$, then show that $\text{Var}(X) \leq (b - a)^2/4$ and that this bound is sharp. (Hints: Reduce to the case $[a, b] = [0, 1]$. If $E(X) = p$, then demonstrate that $\text{Var}(X) \leq p(1 - p)$.)
42. Suppose $g(x)$ is a function such that $g(x) \leq 1$ for all x and $g(x) \leq 0$ for $x \leq c$. Demonstrate the inequality

$$\Pr(X \geq c) \geq E[g(X)] \tag{6.13}$$

for any random variable X [89]. Verify that the polynomial

$$g(x) = \frac{(x - c)(c + 2d - x)}{d^2}$$

with $d > 0$ satisfies the stated conditions leading to inequality (6.13). If X is nonnegative with $E(X) = 1$ and $E(X^2) = \beta$ and $c \in (0, 1)$, then prove that the choice $d = \beta/(1 - c)$ yields

$$\Pr(X \geq c) \geq \frac{(1 - c)^2}{\beta}.$$

Finally, if $E(X^2) = 1$ and $E(X^4) = \beta$, show that

$$\Pr(|X| \geq c) \geq \frac{(1 - c^2)^2}{\beta}.$$

43. Let X be a Poisson random variable with mean λ . Demonstrate that the Chernoff bound

$$\Pr(X \geq c) \leq \inf_{t>0} e^{-ct} E(e^{tX})$$

amounts to

$$\Pr(X \geq c) \leq \frac{(\lambda e)^c}{c^c} e^{-\lambda}$$

for any integer $c > \lambda$. Recall that $\Pr(X = i) = \lambda^i e^{-\lambda}/i!$ for all nonnegative integers i .

44. Use Jensen's inequality to prove the inequality

$$\prod_{k=1}^n x_k^{\alpha_k} + \prod_{k=1}^n y_k^{\alpha_k} \leq \prod_{k=1}^n (x_k + y_k)^{\alpha_k}$$

for positive numbers x_k and y_k and nonnegative numbers α_k with sum $\sum_{k=1}^n \alpha_k = 1$. Prove the inequality

$$\left(1 + \prod_{k=1}^n x_k^{\alpha_k}\right)^{-1} \leq \sum_{k=1}^n \frac{\alpha_k}{1 + x_k}$$

when all $x_k \geq 1$ and the reverse inequality when all $x_k \in (0, 1]$.

45. Let $B_n f(x) = E[f(S_n/n)]$ denote the Bernstein polynomial of degree n approximating $f(x)$ as discussed in Example 6.6.1. Prove that

- (a) $B_n f(x)$ is linear in $f(x)$,
- (b) $B_n f(x) \geq 0$ if $f(x) \geq 0$,
- (c) $B_n f(x) = f(x)$ if $f(x)$ is linear,
- (d) $B_n x(1-x) = \frac{n-1}{n} x(1-x)$,
- (e) $\|B_n f\|_\infty \leq \|f\|_\infty$.

46. Suppose the function $f(x)$ has continuous derivative $f'(x)$. For $\delta > 0$ show that Bernstein's polynomial satisfies the bound

$$\left| E \left[f \left(\frac{S_n}{n} \right) \right] - f(x) \right| \leq \delta \|f'\|_\infty + \frac{\|f\|_\infty}{2n\delta^2}.$$

Conclude from this estimate that $\left\| E \left[f \left(\frac{S_n}{n} \right) \right] - f \right\|_\infty = O(n^{-\frac{1}{3}})$.

47. Let $f(x)$ be a convex function on $[0, 1]$. Prove that the Bernstein polynomial of degree n approximating $f(x)$ is also convex. (Hint: Show that

$$\begin{aligned} \frac{d^2}{dx^2} \mathbb{E} \left[f \left(\frac{S_n}{n} \right) \right] &= n(n-1) \left\{ \mathbb{E} \left[f \left(\frac{S_{n-2} + 2}{n} \right) \right] \right. \\ &\quad \left. - 2 \mathbb{E} \left[f \left(\frac{S_{n-2} + 1}{n} \right) \right] + \mathbb{E} \left[f \left(\frac{S_{n-2}}{n} \right) \right] \right\} \end{aligned}$$

in the notation of Example 6.6.1.)

48. Verify the following special cases

$$\begin{aligned} \sum_{m=1}^n \frac{a_m}{[m(m+1)]^{1/5}} &\leq \left(\sum_{m=1}^n |a_m|^{5/4} \right)^{4/5} \\ \sum_{m=1}^n \frac{a_m}{\sqrt{m}} &\leq \frac{\sqrt{\pi}}{6^{1/4}} \left(\sum_{m=1}^n |a_m|^{4/3} \right)^{3/4} \\ \sum_{m=0}^{\infty} a_m x^m &\leq (1-x^3)^{-1/3} \left(\sum_{m=0}^{\infty} |a_m|^{3/2} \right)^{2/3} \end{aligned}$$

of Hölder's inequality [243]. In the last inequality we take $0 \leq x < 1$.

49. Suppose $1 \leq p < \infty$. For a random variable X with $\mathbb{E}(|X|^p) < \infty$, define the norm $\|X\|_p = \mathbb{E}(X^p)^{1/p}$. Now prove Minkowski's triangle inequality $\|X+Y\|_p \leq \|X\|_p + \|Y\|_p$. (Hint: Apply Hölder's inequality to the right-hand side of

$$\mathbb{E}(|X+Y|^p) \leq \mathbb{E}(|X| \cdot |X+Y|^{p-1}) + \mathbb{E}(|Y| \cdot |X+Y|^{p-1})$$

and rearrange the result.)

50. Suppose X is a random variable satisfying $0 < a \leq X \leq b < \infty$. Use L'Hôpital's rule to prove that the weighted mean $M(p) = \mathbb{E}(X^p)^{1/p}$ is continuous at $p = 0$ if we define $M(0) = e^{\mathbb{E}(\ln X)}$.
51. Suppose the random variable X is bounded below by the positive constant a and above by the positive constant b . Prove Kantorvich's inequality

$$\mathbb{E}(X) \mathbb{E}(X^{-1}) \leq \frac{\mu^2}{\gamma^2}$$

for $\mu = \frac{1}{2}(a+b)$ and $\gamma = \sqrt{ab}$. (Hint: By homogeneity it suffices to consider the case $\gamma = 1$. Apply the arithmetic-geometric mean inequality [243].)

7

Block Relaxation

7.1 Introduction

As a gentle introduction to optimization algorithms, we now consider block relaxation. The more descriptive terms block descent and block ascent suggest either minimization or maximization rather than generic optimization. Regardless of what one terms the strategy, in many problems it pays to update only a subset of the parameters at a time. Block relaxation divides the parameters into disjoint blocks and cycles through the blocks, updating only those parameters within the pertinent block at each stage of a cycle [59]. When each block consists of a single parameter, block relaxation is called cyclic coordinate descent or cyclic coordinate ascent. Block relaxation is best suited to unconstrained problems where the domain of the objective function reduces to a Cartesian product of the subdomains associated with the different blocks. Obviously, exact block updates are a huge advantage. Equality constraints usually present insuperable barriers to coordinate descent and ascent because parameters get locked into position. In some problems it is advantageous to consider overlapping blocks.

The rest of this chapter consists of sequence of examples, most of which are drawn from statistics. Details of statistical inference are downplayed, but familiarity with classical statistics certainly helps in understanding. Block relaxation sometimes converges slowly. In compensation, updates are often very cheap to compute. Judging the performance of optimization algorithms is a complex task. Computational speed is only one factor.

Reliability and ease of implementation can be equally important. In many problems block relaxation is trivial to implement.

7.2 Examples of Block Relaxation

Example 7.2.1 Sinkhorn's Algorithm

Let $\mathbf{M} = (m_{ij})$ be a rectangular matrix with positive entries. Sinkhorn's theorem [237] says that there exist two diagonal matrices \mathbf{A} and \mathbf{B} with positive diagonal entries a_i and b_j such that the matrix \mathbf{AMB} has prescribed row and column sums. Let $r_i > 0$ be the i th row sum and $c_j > 0$ the j th column sum. Because \mathbf{AMB} has entry $a_i m_{ij} b_j$ at the intersection of row i and column j , the constraints are

$$\sum_i a_i m_{ij} b_j = c_j \quad \text{and} \quad \sum_j a_i m_{ij} b_j = r_i.$$

For these constraints to be consistent, we must have

$$\sum_i r_i = \sum_i \sum_j a_i m_{ij} b_j = \sum_j c_j.$$

Given this assumption, we now sketch a method for finding \mathbf{A} and \mathbf{B} . Consider minimizing the smooth function [156]

$$f(\mathbf{A}, \mathbf{B}) = -\sum_i r_i \ln a_i - \sum_j c_j \ln b_j + \sum_i \sum_j a_i m_{ij} b_j.$$

If any a_i or b_j approaches 0, then $f(\mathbf{A}, \mathbf{B})$ tends to ∞ . In view of this fact, the minimum occurs in a region where the parameters a_i and b_j are uniformly bounded below by a positive constant. Within this region, it follows that $a_i m_{ij} b_j$ tends to ∞ if either a_i or b_j tends to ∞ . Hence, the minimum of $f(\mathbf{A}, \mathbf{B})$ exists. At the minimum, Fermat's principle requires

$$\begin{aligned} \frac{\partial}{\partial a_i} f(\mathbf{A}, \mathbf{B}) &= -\frac{r_i}{a_i} + \sum_j m_{ij} b_j = 0 \\ \frac{\partial}{\partial b_j} f(\mathbf{A}, \mathbf{B}) &= -\frac{c_j}{b_j} + \sum_i a_i m_{ij} = 0. \end{aligned}$$

These equations are just a disguised form of Sinkhorn's constraints.

The direct attempt to solve the stationarity equations is almost immediately thwarted. It is much easier to minimize $f(\mathbf{A}, \mathbf{B})$ with respect to \mathbf{A} for \mathbf{B} fixed or vice versa. If we fix \mathbf{B} , then rearranging the first stationarity equation gives

$$a_i = \frac{r_i}{\sum_j m_{ij} b_j}.$$

Similarly, if we fix \mathbf{A} , then rearranging the second stationarity equation yields

$$b_j = \frac{c_j}{\sum_i a_i m_{ij}}.$$

Sinkhorn's block relaxation algorithm [237] alternates the updates of \mathbf{A} and \mathbf{B} . ■

Example 7.2.2 *Poisson Sports Model*

Consider a simplified version of a model proposed by Maher [185] for a sports contest between two teams in which the number of points scored by team i against team j follows a Poisson process with intensity $e^{o_i - d_j}$, where o_i is an “offensive strength” parameter for team i and d_j is a “defensive strength” parameter for team j . (See Sect. 8.9 for a brief description of Poisson processes.) If t_{ij} is the length of time that i plays j and p_{ij} is the number of points that i scores against j , then the corresponding Poisson loglikelihood function is

$$\ell_{ij}(\boldsymbol{\theta}) = p_{ij}(o_i - d_j) + p_{ij} \ln t_{ij} - t_{ij} e^{o_i - d_j} - \ln p_{ij}!, \quad (7.1)$$

where $\boldsymbol{\theta} = (\mathbf{o}, \mathbf{d})$ is the parameter vector. Note that the parameters should satisfy a linear constraint such as $d_1 = 0$ in order for the model be identifiable; otherwise, it is clearly possible to add the same constant to each o_i and d_j without altering the likelihood. We make two simplifying assumptions. First, the outcomes of the different games are independent. Second, each team's point total within a single game is independent of its opponent's point total. The second assumption is more suspect than the first since it implies that a team's offensive and defensive performances are somehow unrelated to one another; nonetheless, the model gives an interesting first approximation to reality. Under these assumptions, the full data loglikelihood is obtained by summing $\ell_{ij}(\boldsymbol{\theta})$ over all pairs (i, j) . Setting the partial derivatives of the loglikelihood equal to zero leads to the equations

$$e^{-d_j} = \frac{\sum_i p_{ij}}{\sum_i t_{ij} e^{o_i}} \quad \text{and} \quad e^{o_i} = \frac{\sum_j p_{ij}}{\sum_j t_{ij} e^{-d_j}}$$

satisfied by the maximum likelihood estimate $(\hat{\mathbf{o}}, \hat{\mathbf{d}})$.

These equations do not admit a closed-form solution, so we turn to block relaxation [59]. If we fix the o_i , then we can solve for the d_j and vice versa in the form

$$d_j = -\ln \left(\frac{\sum_i p_{ij}}{\sum_i t_{ij} e^{o_i}} \right) \quad \text{and} \quad o_i = \ln \left(\frac{\sum_j p_{ij}}{\sum_j t_{ij} e^{-d_j}} \right).$$

Block relaxation consists in alternating the updates of the defensive and offensive parameters with the proviso that d_1 is fixed at 0.

TABLE 7.1. Ranking of all 29 NBA teams on the basis of the 2002–2003 regular season according to their estimated offensive plus defensive strengths. Each team played 82 games

Team	$\hat{o}_i + \hat{d}_i$	Wins	Team	$\hat{o}_i + \hat{d}_i$	Wins
Cleveland	-0.0994	17	Phoenix	0.0166	44
Denver	-0.0845	17	New Orleans	0.0169	47
Toronto	-0.0647	24	Philadelphia	0.0187	48
Miami	-0.0581	25	Houston	0.0205	43
Chicago	-0.0544	30	Minnesota	0.0259	51
Atlanta	-0.0402	35	LA Lakers	0.0277	50
LA Clippers	-0.0355	27	Indiana	0.0296	48
Memphis	-0.0255	28	Utah	0.0299	47
New York	-0.0164	37	Portland	0.0320	50
Washington	-0.0153	37	Detroit	0.0336	50
Boston	-0.0077	44	New Jersey	0.0481	49
Golden State	-0.0051	38	San Antonio	0.0611	60
Orlando	-0.0039	42	Sacramento	0.0686	59
Milwaukee	-0.0027	42	Dallas	0.0804	60
Seattle	0.0039	40			

Table 7.1 summarizes our application of the Poisson sports model to the results of the 2002–2003 regular season of the National Basketball Association. In these data, t_{ij} is measured in minutes. A regular game lasts 48 min, and each overtime period, if necessary, adds 5 min. Thus, team i is expected to score $48e^{\hat{o}_i - \hat{d}_j}$ points against team j when the two teams meet and do not tie. Team i is ranked higher than team j if $\hat{o}_i - \hat{d}_j > \hat{o}_j - \hat{d}_i$, which is equivalent to the condition $\hat{o}_i + \hat{d}_i > \hat{o}_j + \hat{d}_j$.

It is worth emphasizing some of the virtues of the model. First, the ranking of the 29 NBA teams on the basis of the estimated sums $\hat{o}_i + \hat{d}_i$ for the 2002–2003 regular season is not perfectly consistent with their cumulative wins; strength of schedule and margins of victory are reflected in the model. Second, the model gives the point-spread function for a particular game as the difference of two independent Poisson random variables. Third, one can easily amend the model to rank individual players rather than teams by assigning to each player an offensive and defensive intensity parameter. If each game is divided into time segments punctuated by substitutions, then the block relaxation algorithm can be adapted to estimate the assigned player intensities. This might provide a rational basis for salary negotiations that takes into account subtle differences between players not reflected in traditional sports statistics. ■

Example 7.2.3 *K-Means Clustering*

In k-means clustering we must divide n points $\mathbf{x}_1, \dots, \mathbf{x}_n$ in \mathbb{R}^m into k clusters. Each cluster C_j is characterized by a cluster center $\boldsymbol{\mu}_j$. The best

clustering of the points minimizes the criterion

$$f(\boldsymbol{\mu}, C) = \sum_{j=1}^k \sum_{\mathbf{x}_i \in C_j} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|^2,$$

where $\boldsymbol{\mu}$ is the matrix whose columns are the $\boldsymbol{\mu}_j$ and C is the collection of clusters. Because this mixed continuous-discrete optimization problem has no obvious analytic solution, block relaxation is attractive. If we hold the clusters fixed, then it is clear from Example 6.5.7 that we should set

$$\boldsymbol{\mu}_j = \frac{1}{|C_j|} \sum_{\mathbf{x}_i \in C_j} \mathbf{x}_i.$$

Similarly, it is clear that if we hold the cluster centers fixed, then we should assign point \mathbf{x}_i to the cluster C_j minimizing $\|\mathbf{x}_i - \boldsymbol{\mu}_j\|$. Block relaxation, known as Lloyd's algorithm in this context, alternates cluster center redefinition and cluster membership reassignment. It is simple and effective. The initial cluster centers can be chosen randomly from the n data points. The evidence suggests that this should be done in a biased manner that spreads the centers out [5]. Changing the objective function to

$$g(\boldsymbol{\mu}, C) = \sum_{j=1}^k \sum_{\mathbf{x}_i \in C_j} \|\mathbf{x}_i - \boldsymbol{\mu}_j\|_1$$

makes it more resistant to outliers. The recentering step is now solved by replacing means by medians in each coordinate. This takes a little more computation but is usually worth the effort. ■

Example 7.2.4 Canonical Correlations

Consider a random vector \mathbf{Z} partitioned into a subvector \mathbf{X} of predictors and a subvector \mathbf{Y} of responses. (See Sect. 9.7 for a brief discussion of random vectors, expectations, and variances.) The most elementary form of canonical correlation analysis seeks two linear combinations $\mathbf{a}^* \mathbf{X}$ and $\mathbf{b}^* \mathbf{Y}$ that are maximally correlated [187]. If we partition the variance matrix of \mathbf{Z} into blocks

$$\text{Var}(\mathbf{Z}) = \begin{pmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{pmatrix}$$

consistent with \mathbf{X} and \mathbf{Y} , then the two linear combinations maximize the covariance $\mathbf{a}^* \boldsymbol{\Sigma}_{12} \mathbf{b}$ subject to the variance constraints

$$\mathbf{a}^* \boldsymbol{\Sigma}_{11} \mathbf{a} = \mathbf{b}^* \boldsymbol{\Sigma}_{22} \mathbf{b} = 1.$$

This constrained maximization problem is an ideal candidate for block relaxation. Problems 8 and 9 relate the best vectors \mathbf{a} and \mathbf{b} to the singular value decomposition of a matrix.

TABLE 7.2. Iterates in canonical correlation estimation

n	a_{n1}	a_{n2}	b_{n1}	b_{n2}
0	1.000000	1.000000	1.000000	1.000000
1	0.553047	0.520658	0.504588	0.538164
2	0.552159	0.521554	0.504509	0.538242
3	0.552155	0.521558	0.504509	0.538242
4	0.552155	0.521558	0.504509	0.538242

For fixed \mathbf{b} we can easily find the best \mathbf{a} . Introduce the Lagrangian

$$\mathcal{L}(\mathbf{a}) = \mathbf{a}^* \Sigma_{12} \mathbf{b} - \frac{\lambda}{2} (\mathbf{a}^* \Sigma_{11} \mathbf{a} - 1),$$

and equate its gradient

$$\nabla \mathcal{L}(\mathbf{a}) = \Sigma_{12} \mathbf{b} - \lambda \Sigma_{11} \mathbf{a}$$

to $\mathbf{0}$. This gives the maximum point

$$\mathbf{a} = \frac{1}{\lambda} \Sigma_{11}^{-1} \Sigma_{12} \mathbf{b},$$

assuming the submatrix Σ_{11} is positive definite. Inserting this value into the constraint $\mathbf{a}^* \Sigma_{11} \mathbf{a} = 1$ allows us to solve for the Lagrange multiplier λ and hence pin down \mathbf{a} as

$$\mathbf{a} = \frac{1}{\sqrt{\mathbf{b}^* \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} \mathbf{b}}} \Sigma_{11}^{-1} \Sigma_{12} \mathbf{b}.$$

Because the second differential $d^2 \mathcal{L} = -\lambda \Sigma_{11}$ is negative definite, \mathbf{a} represents the maximum. Likewise, fixing \mathbf{a} and optimizing over \mathbf{b} gives the update

$$\mathbf{b} = \frac{1}{\sqrt{\mathbf{a}^* \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}}} \Sigma_{22}^{-1} \Sigma_{21} \mathbf{a}.$$

As a toy example consider the correlation matrix

$$\text{Var}(\mathbf{Z}) = \begin{pmatrix} 1 & 0.7346 & 0.7108 & 0.7040 \\ 0.7346 & 1 & 0.6932 & 0.7086 \\ 0.7108 & 0.6932 & 1 & 0.8392 \\ 0.7040 & 0.7086 & 0.8392 & 1 \end{pmatrix}$$

with unit variances on its diagonal. Table 7.2 shows the first few iterates of block relaxation starting from $\mathbf{a} = \mathbf{b} = \mathbf{1}$. Convergence is exceptionally quick; more complex examples exhibit slower convergence. ■

Example 7.2.5 *Iterative Proportional Fitting*

Our next example of block relaxation is taken from the contingency table literature [14, 86]. Consider a three-way contingency table with two-way interactions. If the three factors are indexed by i , j , and k and have r , s , and t levels, respectively, then a loglinear model for the observed count data y_{ijk} is defined by an exponentially parameterized mean

$$\mu_{ijk} = e^{\lambda + \lambda_i^1 + \lambda_j^2 + \lambda_k^3 + \lambda_{ij}^{12} + \lambda_{ik}^{13} + \lambda_{jk}^{23}}$$

for each cell ijk . To ensure that all parameters are identifiable, we make the usual assumption that a parameter set summed over one of its indices yields 0. For instance, $\lambda_i^1 = \sum_i \lambda_i^1 = 0$ and $\lambda_{ij}^{12} = \sum_j \lambda_{ij}^{12} = 0$. The overall effect λ is permitted to be nonzero.

If we postulate independent Poisson distributions for the random variables Y_{ijk} underlying the observed values y_{ijk} , then the loglikelihood is

$$L = \sum_i \sum_j \sum_k (y_{ijk} \ln \mu_{ijk} - \mu_{ijk}). \quad (7.2)$$

Maximizing L with respect to λ can be accomplished by setting

$$\frac{\partial}{\partial \lambda} L = \sum_i \sum_j \sum_k (y_{ijk} - \mu_{ijk}) = 0.$$

This tells us that whatever the other parameters are, λ should be adjusted so that $\mu_{...} = y_{...} = m$ is the total sample size. (Here again the dot convention signifies summation over a lost index.) In other words, if $\mu_{ijk} = e^{\lambda} \omega_{ijk}$, then λ is chosen so that $e^{\lambda} = m/\omega_{...}$. With this proviso, the loglikelihood becomes

$$\begin{aligned} L &= \sum_i \sum_j \sum_k y_{ijk} \ln \frac{m \omega_{ijk}}{\omega_{...}} - m \\ &= \sum_i \sum_j \sum_k y_{ijk} \ln \frac{\omega_{ijk}}{\omega_{...}} + m \ln m - m, \end{aligned}$$

which is up to an irrelevant constant just the loglikelihood of a multinomial distribution with probability $\omega_{ijk}/\omega_{...}$ attached to cell ijk . Thus, for purposes of maximum likelihood estimation, we might as well stick with the Poisson sampling model.

Unfortunately, no closed-form solution to the Poisson likelihood equations exists satisfying the complicated linear constraints. The resolution of this dilemma lies in refusing to update all of the parameters simultaneously. Suppose that we consider only the parameters λ , λ_i^1 , λ_j^2 , and λ_{ij}^{12} pertinent to the first two factors. If in equation (7.2) we let

$$\begin{aligned} \mu_{ij} &= e^{\lambda + \lambda_i^1 + \lambda_j^2 + \lambda_{ij}^{12}} \\ \alpha_{ijk} &= e^{\lambda_k^3 + \lambda_{ik}^{13} + \lambda_{jk}^{23}}, \end{aligned}$$

then setting

$$\begin{aligned} \frac{\partial}{\partial \lambda_{ij}^{12}} L &= \sum_k (y_{ijk} - \mu_{ijk}) \\ &= y_{ij.} - \mu_{ij.} \\ &= y_{ij.} - \mu_{ij} \alpha_{ij.} \\ &= 0 \end{aligned}$$

leads to $\mu_{ij} = y_{ij.}/\alpha_{ij.}$. The constraint $\sum_k (y_{ijk} - \mu_{ijk}) = 0$ implies that the other partial derivatives

$$\begin{aligned} \frac{\partial}{\partial \lambda} L &= y_{...} - \mu_{...} \\ \frac{\partial}{\partial \lambda_i^1} L &= y_{i..} - \mu_{i..} \\ \frac{\partial}{\partial \lambda_j^2} L &= y_{.j.} - \mu_{.j.} \end{aligned}$$

vanish as well. This stationary point of the loglikelihood is also a stationary point of the Lagrangian with all Lagrange multipliers equal to 0.

Of course, we still must nail down λ , λ_i^1 , λ_j^2 , and λ_{ij}^{12} . In view of the definition of μ_{ij} , the choice

$$\lambda_{ij}^{12} = \ln \left(\frac{y_{ij.}}{\alpha_{ij.}} \right) - \lambda - \lambda_i^1 - \lambda_j^2$$

guarantees that $\mu_{ij} = y_{ij.}/\alpha_{ij.}$. One can check that the further choices

$$\begin{aligned} \lambda &= \frac{1}{rs} \sum_i \sum_j \ln \mu_{ij} \\ \lambda_i^1 &= \frac{1}{s} \sum_j \ln \mu_{ij} - \lambda \\ \lambda_j^2 &= \frac{1}{r} \sum_i \ln \mu_{ij} - \lambda \end{aligned}$$

satisfy the relevant equality constraints $\lambda^1 = 0$, $\lambda^2 = 0$, $\lambda_j^{12} = 0$, and $\lambda_i^{12} = 0$. The identity $\mu_{ij} = y_{ij.}/\alpha_{ij.}$ is crucial in this regard.

At the second stage, the parameter set $\{\lambda, \lambda_i^1, \lambda_k^3, \lambda_{ik}^{13}\}$ is updated, holding the remaining parameters fixed. At the third stage, the parameter set $\{\lambda, \lambda_j^2, \lambda_k^3, \lambda_{jk}^{23}\}$ is updated, holding the remaining parameters fixed. These three successive stages constitute one iteration of the iterative proportional fitting algorithm. Each stage either leaves all parameters unchanged or increases the loglikelihood. In this example, the parameter blocks are not disjoint. ■

Example 7.2.6 *Matrix Factorization by Alternating Least Squares*

Least squares, the most venerable of the statistical fitting procedures, was initiated by Gauss and Legendre. As explained in Example 1.3.3, the basic setup involves n independent responses that individually take the form

$$y_i = \sum_{j=1}^p x_{ij}\theta_j + u_i. \quad (7.3)$$

Here y_i depends linearly on the unknown regression coefficients θ_j through the known predictors x_{ij} . The error u_i is assumed to be normally distributed with mean 0 and variance σ^2 . If we collect the y_i into a $n \times 1$ response vector \mathbf{y} , the x_{ij} into a $n \times p$ design matrix \mathbf{X} , the θ_j into a $p \times 1$ parameter vector $\boldsymbol{\theta}$, and the u_j into a $n \times 1$ error vector \mathbf{u} , then the linear regression model can be rewritten in vector notation as $\mathbf{y} = \mathbf{X}\boldsymbol{\theta} + \mathbf{u}$. Provided the design matrix \mathbf{X} has full rank, the least squares estimate of $\boldsymbol{\theta}$ is $\hat{\boldsymbol{\theta}} = (\mathbf{X}^* \mathbf{X})^{-1} \mathbf{X}^* \mathbf{y}$. Example 5.2.6 shows how to amend this solution to take into account affine constraints. In weighted least squares, one minimizes the criterion

$$f(\boldsymbol{\theta}) = \sum_{i=1}^n c_i \left(y_i - \sum_{j=1}^p x_{ij}\theta_j \right)^2,$$

where the c_i are positive weights. This reduces to ordinary least squares if one substitutes $\sqrt{c_i}y_i$ for y_i and $\sqrt{c_i}x_{ij}$ for x_{ij} . It is clear that any method for solving an ordinary least squares problem can be immediately adapted to solving a weighted least squares problem.

The history of alternating least squares is summarized by Gifi [104]. Very early on Kruskal [158] applied the method to factorial ANOVA. Here we briefly survey its use in nonnegative matrix factorization [174, 175]. Suppose \mathbf{U} is a $n \times q$ matrix whose columns $\mathbf{u}_1, \dots, \mathbf{u}_q$ represent data vectors. In many applications one wants to explain the data by postulating a reduced number of prototypes $\mathbf{v}_1, \dots, \mathbf{v}_p$ and writing

$$\mathbf{u}_j \approx \sum_{k=1}^p \mathbf{v}_k w_{kj}$$

for certain nonnegative weights w_{kj} . The matrix $\mathbf{W} = (w_{kj})$ is $p \times q$. If p is small compared to q , then the representation $\mathbf{U} \approx \mathbf{V}\mathbf{W}$ compresses the data for easier storage and retrieval. Depending on the circumstances, further constraints may be advisable [72]. For instance, if the entries of \mathbf{U} are nonnegative, then it is often reasonable to demand that the entries of \mathbf{V} be nonnegative as well. If we want each \mathbf{u}_j to equal a convex combination of the prototypes, then constraining the column sums of \mathbf{W} to equal 1 is indicated.

One way of estimating \mathbf{V} and \mathbf{W} is to minimize the objective function

$$\|\mathbf{U} - \mathbf{V}\mathbf{W}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^q \left(u_{ij} - \sum_{k=1}^p v_{ik} w_{kj} \right)^2.$$

No explicit solution is known, but alternating least squares offers an iterative attack. If \mathbf{W} is fixed, then we can update the i th row of \mathbf{V} by minimizing the sum of squares

$$\sum_{j=1}^q \left(u_{ij} - \sum_{k=1}^p v_{ik} w_{kj} \right)^2.$$

Similarly, if \mathbf{V} is fixed, then we can update the j th column of \mathbf{W} by minimizing the sum of squares

$$\sum_{i=1}^n \left(u_{ij} - \sum_{k=1}^p v_{ik} w_{kj} \right)^2.$$

In either case we are faced with solving a sequence of least squares problems. The introduction of nonnegativity constraints and convexity constraints complicates matters. Problem 12 suggests coordinate descent methods for solving these two constrained least squares problems. Coordinate descent is trivial to implement but potentially very slow. We will revisit nonnegative least squares later from a different perspective. ■

7.3 Problems

1. Program and test any one of the six examples in this chapter.
2. Demonstrate that cyclic coordinate descent either diverges or converges to a saddle point of the function $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ defined by

$$f(\mathbf{x}) = (x_1 - x_2)^2 - 2x_1x_2.$$

This function of de Leeuw [59] has no minimum.

3. Consider the function $f(\mathbf{x}) = (x_1^2 + x_2^2)^{-1} + \ln(x_1^2 + x_2^2)$ for $\mathbf{x} \neq \mathbf{0}$. Explicitly find the minimum value of $f(\mathbf{x})$. Specify the coordinate descent algorithm for finding the minimum. Note any ambiguities in the implementation of coordinate descent, and describe the possible cluster points of the algorithm as a function of the initial point. (Hint: Coordinate descent, properly defined, converges in a finite number of iterations.)

4. Implement cyclic coordinate descent to minimize the posynomial

$$f(\mathbf{x}) = \frac{1}{x_1^3} + \frac{3}{x_1 x_2^2} + x_1 x_2$$

over the region $\{\mathbf{x} \in \mathbb{R}^2 : x_1 > 0, x_2 > 0\}$. Derive the updates

$$x_{n+1,1} = \sqrt{\frac{\frac{3}{x_{n2}^2} + \sqrt{\frac{9}{x_{n2}^4} + 12x_{n2}}}{2x_{n2}}}$$

$$x_{n+1,2} = \sqrt[3]{\frac{6}{x_{n+1,1}^2}}.$$

Compare your numerical results to those displayed in Table 8.2.

5. In Sinkhorn's theorem, suppose the matrix \mathbf{M} is square. Show that some entries of \mathbf{M} can be 0 as long as some positive power \mathbf{M}^P of \mathbf{M} has all entries positive.
6. Consider cluster analysis on the real line. Show that Lloyd's algorithm cannot improve on the initial partition $\pi_1 = \{0, 2\}$ and $\pi_2 = \{3, 5\}$ despite the fact that the partition $\pi_1 = \{0\}$ and $\pi_2 = \{2, 3, 5\}$ is better. Also demonstrate that Lloyd's algorithm transforms the initial partition $\pi_1 = \{-7, -5\}$, $\pi_2 = \{-4, 4\}$, and $\pi_3 = \{5, 7\}$ into the partition $\pi_1 = \{-7, -5, -4\}$, $\pi_2 = \emptyset$, and $\pi_3 = \{4, 5, 7\}$.
7. Let $\mathbf{M} = (m_{ij})$ be a nontrivial $m \times n$ matrix. The dominant part of the singular value decomposition (svd) of \mathbf{M} is an outer product matrix $\lambda \mathbf{u} \mathbf{v}^*$ with $\lambda > 0$ and \mathbf{u} and \mathbf{v} unit vectors. This outer product minimizes

$$\|\mathbf{M} - \lambda \mathbf{u} \mathbf{v}^*\|_F^2 = \sum_i \sum_j (m_{ij} - \lambda u_i v_j)^2.$$

One can use alternating least squares to find $\lambda \mathbf{u} \mathbf{v}^*$ [101]. In the first step of the algorithm, one fixes \mathbf{v} and estimates $\mathbf{w} = \lambda \mathbf{u}$ by least squares. Show that \mathbf{w} has components $w_i = \sum_j m_{ij} v_j$. Once \mathbf{w} is available, we set $\lambda = \|\mathbf{w}\|$ and $\mathbf{u} = \|\mathbf{w}\|^{-1} \mathbf{w}$. What are the corresponding updates for \mathbf{v} and λ when you fix \mathbf{u} ? To find the next outer product in the svd, form the deflated matrix $\mathbf{M} - \lambda \mathbf{u} \mathbf{v}^*$ and repeat the process. Program and test this algorithm.

8. Continuing Problem 7, prove that minimizing $\|\mathbf{M} - \lambda \mathbf{u} \mathbf{v}^*\|_F^2$ subject to the constraints $\|\mathbf{u}\| = \|\mathbf{v}\| = 1$ is equivalent to maximizing $\mathbf{u}^* \mathbf{M} \mathbf{v}$ subject to the same constraints. The solution vectors \mathbf{u} and \mathbf{v} are called singular vectors. The corresponding scalar λ is the singular value. Problem 7 of Chap. 2 relates λ to the spectral norm of \mathbf{M} .

9. In Example 7.2.4, make the linear change of variables $\mathbf{c} = \Sigma_{11}^{1/2} \mathbf{a}$ and $\mathbf{d} = \Sigma_{22}^{1/2} \mathbf{b}$. In the new variables show that one must maximize $\mathbf{c}^* \mathbf{\Omega} \mathbf{d}$ subject to $\|\mathbf{c}\| = \|\mathbf{d}\| = 1$, where $\mathbf{\Omega} = \Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}$. In the language of Problems 7 and 8, the first singular vectors $\mathbf{c} = \mathbf{u}$ and $\mathbf{d} = \mathbf{v}$ of the svd of $\mathbf{\Omega}$ solve the transformed problem. Obviously, the vectors $\mathbf{b} = \Sigma_{11}^{-1/2} \mathbf{u}$ and $\mathbf{a} = \Sigma_{22}^{-1/2} \mathbf{v}$ solve the original problem. The advantage of this approach is that one can now define higher-order canonical correlations from the remaining singular vectors of the svd.
10. Suppose \mathbf{A} is a symmetric matrix and \mathbf{B} is a positive definite matrix of the same dimension. Formulate cyclic coordinate descent and ascent algorithms for minimizing and maximizing the Rayleigh quotient

$$R(\mathbf{x}) = \frac{\mathbf{x}^* \mathbf{A} \mathbf{x}}{\mathbf{x}^* \mathbf{B} \mathbf{x}} \tag{7.4}$$

over the set $\mathbf{x} \neq \mathbf{0}$. Program and test this algorithm.

11. Continuing Problem 10, demonstrate that the maximum and minimum values of the Rayleigh quotient (7.4) coincide with the maximum and minimum eigenvalues of the matrix $\mathbf{B}^{-1} \mathbf{A}$.
12. For a positive definite matrix \mathbf{A} , consider minimizing the quadratic function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ subject to the constraints $x_i \geq 0$ for all i . Show that the cyclic coordinate descent updates are

$$\hat{x}_i = \max \left\{ 0, x_i - a_{ii}^{-1} \left[\sum_j a_{ij} x_j + b_i \right] \right\}.$$

If we impose the additional constraint $\sum_i x_i = 1$, the problem is harder. One line of attack is to minimize the penalized function

$$f_\mu(\mathbf{x}) = f(\mathbf{x}) + \frac{\mu}{2} \left(\sum_i x_i - 1 \right)^2$$

for a large positive constant μ . The theory in Chap. 13 shows that the minimum of $f_\mu(\mathbf{x})$ tends to the constrained minimum of $f(\mathbf{x})$ as μ tends to ∞ . Accepting this result, demonstrate that cyclic coordinate descent for $f_\mu(\mathbf{x})$ has updates

$$\hat{x}_i = \max \left\{ 0, x_i - (a_{ii} + \mu)^{-1} \left[\sum_j a_{ij} x_j + b_i + \mu \left(\sum_j x_j - 1 \right) \right] \right\}.$$

Program this second algorithm and test it for the choices

$$\mathbf{A} = \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix}, \quad \mathbf{b} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}.$$

Start with $\mu = 1$ and double it every time you update the full vector \mathbf{x} . Do the iterates converge to the minimum of $f(\mathbf{x})$ subject to all constraints?

TABLE 7.3. Coronary disease data

Disease status	Cholesterol level	Blood pressure				Total
		1	2	3	4	
Coronary	1	2	3	3	4	12
	2	3	2	1	3	9
	3	8	11	6	6	31
	4	7	12	11	11	41
Total		20	28	21	24	93
No coronary	1	117	121	47	22	307
	2	85	98	43	20	246
	3	119	209	68	43	439
	4	67	99	46	33	245
Total		388	527	204	118	1,237

13. Program and test a k-medians clustering algorithm and concoct an example where it differs from k-means clustering.
14. In fitting splines to data, the problem arises of minimizing the criterion $\|\mathbf{y} - \mathbf{U}\boldsymbol{\alpha} - \mathbf{V}\boldsymbol{\beta}\|^2 + \lambda\boldsymbol{\beta}^* \mathbf{W}\boldsymbol{\beta}$ with respect $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ for $\lambda \geq 0$ and \mathbf{W} positive semidefinite [166]. Derive the block descent updates

$$\begin{aligned}\boldsymbol{\alpha} &= (\mathbf{U}^* \mathbf{U})^{-1} \mathbf{U}^* (\mathbf{y} - \mathbf{V}\boldsymbol{\beta}) \\ \boldsymbol{\beta} &= (\mathbf{V}^* \mathbf{V} + \lambda \mathbf{W})^{-1} \mathbf{V}^* (\mathbf{y} - \mathbf{U}\boldsymbol{\alpha}).\end{aligned}$$

15. Consider the coronary disease data [86, 159] displayed in the three-way contingency Table 7.3. Using iterative proportional fitting, find the maximum likelihood estimates for the loglinear model with first-order interactions. Perform a chi-square test to decide whether this model fits the data better than the model postulating independence of the three factors.
16. As noted in the text, the loglinear model for categorical data can be interpreted as assuming independent Poisson distributions for the various categories with category i having mean $\mu_i(\boldsymbol{\theta}) = e^{\mathbf{l}_i^* \boldsymbol{\theta}}$, where \mathbf{l}_i is a vector whose entries are 0's or 1's. Calculate the observed information $-d^2 L(\boldsymbol{\theta}) = \sum_i e^{\mathbf{l}_i^* \boldsymbol{\theta}} \mathbf{l}_i \mathbf{l}_i^*$ in this circumstance, and deduce that it is positive semidefinite. In the presence of affine constraints $\mathbf{V}\boldsymbol{\theta} = \mathbf{d}$ on $\boldsymbol{\theta}$, show that any maximum likelihood estimate of $\boldsymbol{\theta}$ is necessarily unique provided the null space (kernel) of \mathbf{V} is contained in the linear span of the \mathbf{l}_i .

8

The MM Algorithm

8.1 Introduction

Most practical optimization problems defy exact solution. In the current chapter we discuss an optimization method that relies heavily on convexity arguments and is particularly useful in high-dimensional problems such as image reconstruction [171]. This iterative method is called the MM algorithm. One of the virtues of this acronym is that it does double duty. In minimization problems, the first M of MM stands for majorize and the second M for minimize. In maximization problems, the first M stands for minorize and the second M for maximize. When it is successful, the MM algorithm substitutes a simple optimization problem for a difficult optimization problem. Simplicity can be attained by: (a) separating the variables of an optimization problem, (b) avoiding large matrix inversions, (c) linearizing an optimization problem, (d) restoring symmetry, (e) dealing with equality and inequality constraints gracefully, and (f) turning a non-differentiable problem into a smooth problem. In simplifying the original problem, we must pay the price of iteration or iteration with a slower rate of convergence.

Statisticians have vigorously developed a special case of the MM algorithm called the EM algorithm, which revolves around notions of missing data [65, 166, 191]. We present the EM algorithm in the next chapter. We prefer to present the MM algorithm first because of its greater generality, its more obvious connection to convexity, and its weaker reliance on difficult statistical principles.

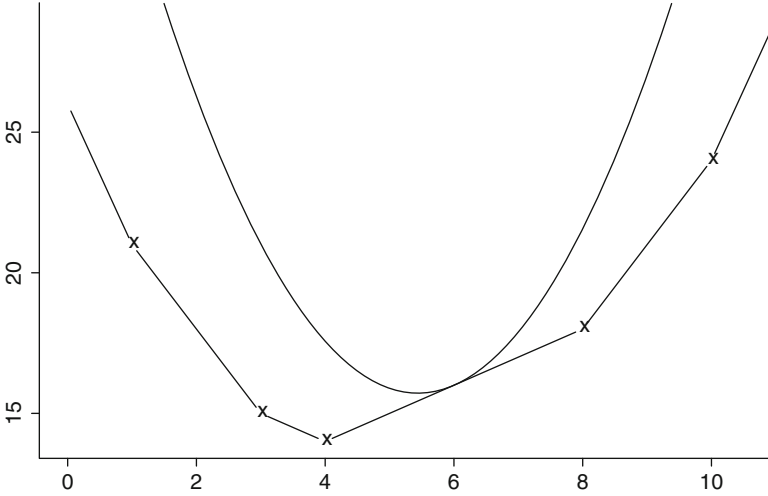


FIGURE 8.1. A quadratic majorizing function for the piecewise linear function $f(x) = |x - 1| + |x - 3| + |x - 4| + |x - 8| + |x - 10|$ at the point $x_m = 6$

8.2 Philosophy of the MM Algorithm

A function $g(\mathbf{x} \mid \mathbf{x}_m)$ is said to majorize a function $f(\mathbf{x})$ at \mathbf{x}_m provided

$$f(\mathbf{x}_m) = g(\mathbf{x}_m \mid \mathbf{x}_m) \tag{8.1}$$

$$f(\mathbf{x}) \leq g(\mathbf{x} \mid \mathbf{x}_m), \quad \mathbf{x} \neq \mathbf{x}_m. \tag{8.2}$$

In other words, the surface $\mathbf{x} \mapsto g(\mathbf{x} \mid \mathbf{x}_m)$ lies above the surface $f(\mathbf{x})$ and is tangent to it at the point $\mathbf{x} = \mathbf{x}_m$. Here \mathbf{x}_m represents the current iterate in a search of the surface $f(\mathbf{x})$. Figure 8.1 provides a simple one-dimensional example.

In the minimization version of the MM algorithm, we minimize the surrogate majorizing function $g(\mathbf{x} \mid \mathbf{x}_m)$ rather than the actual function $f(\mathbf{x})$. If \mathbf{x}_{m+1} denotes the minimum of the surrogate $g(\mathbf{x} \mid \mathbf{x}_m)$, then we can show that the MM procedure forces $f(\mathbf{x})$ downhill. Indeed, the inequalities

$$f(\mathbf{x}_{m+1}) \leq g(\mathbf{x}_{m+1} \mid \mathbf{x}_m) \leq g(\mathbf{x}_m \mid \mathbf{x}_m) = f(\mathbf{x}_m) \tag{8.3}$$

follow directly from the definition of \mathbf{x}_{m+1} and the majorization conditions (8.1) and (8.2). The descent property (8.3) lends the MM algorithm remarkable numerical stability. Strictly speaking, it depends only on decreasing $g(\mathbf{x} \mid \mathbf{x}_m)$, not on minimizing $g(\mathbf{x} \mid \mathbf{x}_m)$. This fact has practical consequences when the minimum of $g(\mathbf{x} \mid \mathbf{x}_m)$ cannot be found exactly. When $f(\mathbf{x})$ is strictly convex, one can show with a few additional mild hypotheses that the iterates \mathbf{x}_m converge to the global minimum of $f(\mathbf{x})$ regardless of the initial point \mathbf{x}_0 .

If $g(\mathbf{x} \mid \mathbf{x}_m)$ majorizes $f(\mathbf{x})$ at an interior point \mathbf{x}_m of the domain of $f(\mathbf{x})$, then \mathbf{x}_m is a stationary point of the difference $g(\mathbf{x} \mid \mathbf{x}_m) - f(\mathbf{x})$, and the gradient identity

$$\nabla g(\mathbf{x}_m \mid \mathbf{x}_m) = \nabla f(\mathbf{x}_m) \quad (8.4)$$

holds. Furthermore, the second differential $d^2g(\mathbf{x}_m \mid \mathbf{x}_m) - d^2f(\mathbf{x}_m)$ is positive semidefinite. Problem 1 makes the point that the majorization relation between functions is closed under the formation of sums, nonnegative products, limits, and composition with an increasing function. These rules permit us to work piecemeal in simplifying complicated objective functions. With obvious changes, the MM algorithm also applies to maximization rather than to minimization. To maximize a function $f(\mathbf{x})$, we minorize it by a surrogate function $g(\mathbf{x} \mid \mathbf{x}_m)$ and maximize $g(\mathbf{x} \mid \mathbf{x}_m)$ to produce the next iterate \mathbf{x}_{m+1} .

The reader might well object that the MM algorithm is not so much an algorithm as a vague philosophy for deriving an algorithm. The same objection applies to the EM algorithm. As we proceed through the current chapter, we hope the various examples will convince the reader of the value of a unifying principle and a framework for attacking concrete problems. The strong connection of the MM algorithm to convexity and inequalities has the natural pedagogical advantage of building on the material presented in previous chapters.

8.3 Majorization and Minorization

We will feature five methods for constructing majorizing functions. Two of these simply adapt Jensen's inequality

$$f\left(\sum_i \alpha_i t_i\right) \leq \sum_i \alpha_i f(t_i)$$

defining a convex function $f(t)$. It is easy to identify convex functions on the real line, so the first method composes such a function with a linear function $\mathbf{c}^* \mathbf{x}$ to create a new convex function of the vector \mathbf{x} . Invoking the definition of convexity with $\alpha_i = c_i y_i / \mathbf{c}^* \mathbf{y}$ and $t_i = \mathbf{c}^* \mathbf{y}_i x_i / y_i$ then yields

$$f(\mathbf{c}^* \mathbf{x}) \leq \sum_i \frac{c_i y_i}{\mathbf{c}^* \mathbf{y}} f\left(\frac{\mathbf{c}^* \mathbf{y}}{y_i} x_i\right) = g(\mathbf{x} \mid \mathbf{y}), \quad (8.5)$$

provided all of the components of the vectors \mathbf{c} , \mathbf{x} , and \mathbf{y} are positive. The surrogate function $g(\mathbf{x} \mid \mathbf{y})$ equals $f(\mathbf{c}^* \mathbf{y})$ when $\mathbf{x} = \mathbf{y}$. One of the virtues of applying inequality (8.5) in defining a surrogate function is that it separates parameters in the surrogate function. This feature is critically important in high-dimensional problems because it reduces optimization

over \mathbf{x} to a sequence of one-dimensional optimizations over each component x_i . The argument establishing inequality (8.5) is equally valid if we replace the parameter vector \mathbf{x} throughout by a vector-valued function $h(\mathbf{x})$ of \mathbf{x} . The genetics problem in the next section illustrates this variant of the technique.

To relax the positivity restrictions on the vectors \mathbf{c} , \mathbf{x} , and \mathbf{y} , De Pierro [67] suggested in a medical imaging context the alternative majorization

$$f(\mathbf{c}^* \mathbf{x}) \leq \sum_i \alpha_i f \left\{ \frac{c_i}{\alpha_i} (x_i - y_i) + \mathbf{c}^* \mathbf{y} \right\} = g(\mathbf{x} | \mathbf{y}) \quad (8.6)$$

for a convex function $f(t)$. Here all $\alpha_i \geq 0$, $\sum_i \alpha_i = 1$, and $\alpha_i > 0$ whenever $c_i \neq 0$. In practice, we must somehow tailor the α_i to the problem at hand. Among the obvious candidates for the α_i are

$$\alpha_i = \frac{|c_i|^p}{\sum_j |c_j|^p}$$

for $p \geq 0$. When $p = 0$, we interpret α_i as 0 if $c_i = 0$ and as $1/q$ if c_i is one among q nonzero coefficients.

Our third method involves the linear majorization

$$f(\mathbf{x}) \leq f(\mathbf{y}) + df(\mathbf{y})(\mathbf{x} - \mathbf{y}) = g(\mathbf{x} | \mathbf{y}) \quad (8.7)$$

satisfied by any concave function $f(\mathbf{x})$. Once again we can replace the argument \mathbf{x} by a vector-valued function $h(\mathbf{x})$.

Our fourth method applies to functions $f(\mathbf{x})$ with bounded curvature [16, 59]. Assuming that $f(\mathbf{x})$ is twice differentiable, we look for a matrix \mathbf{B} satisfying $\mathbf{B} \succeq d^2 f(\mathbf{x})$ and $\mathbf{B} \succ \mathbf{0}$ in the sense that $\mathbf{B} - d^2 f(\mathbf{x})$ is positive semidefinite for all \mathbf{x} and \mathbf{B} is positive definite. The quadratic bound principle then amounts to the majorization

$$\begin{aligned} f(\mathbf{x}) &= f(\mathbf{y}) + df(\mathbf{y})(\mathbf{x} - \mathbf{y}) \\ &\quad + (\mathbf{x} - \mathbf{y})^* \int_0^1 d^2 f[\mathbf{y} + t(\mathbf{x} - \mathbf{y})](1-t) dt (\mathbf{x} - \mathbf{y}) \\ &\leq f(\mathbf{y}) + df(\mathbf{y})(\mathbf{x} - \mathbf{y}) + \frac{1}{2} (\mathbf{x} - \mathbf{y})^* \mathbf{B} (\mathbf{x} - \mathbf{y}) \\ &= g(\mathbf{x} | \mathbf{y}). \end{aligned} \quad (8.8)$$

Our fifth and final method exploits the generalized arithmetic-geometric mean inequality (6.5) of Chap. 6. With this result in mind, Problem 9 asks the reader to prove the majorization

$$\prod_{i=1}^n x_i^{\alpha_i} \leq \left(\prod_{i=1}^n y_i^{\alpha_i} \right) \sum_{i=1}^n \frac{\alpha_i}{\alpha} \left(\frac{x_i}{y_i} \right)^\alpha = g(\mathbf{x} | \mathbf{y}) \quad (8.9)$$

for positive numbers x_i , y_i , and α_i and sum $\alpha = \sum_{i=1}^n \alpha_i$. Inequality (8.9) is the key to separating parameters with posynomials. We will use it in sketching an MM algorithm for unconstrained geometric programming.

Any of the first four majorizations can be turned into minorizations by interchanging the adjectives convex and concave and positive definite and negative definite, respectively. Of course, there is an art to applying these methods just as there is an art to applying any mathematical principle. The methods hardly exhaust the possibilities for majorization and minorization. The problems at the end of the chapter sketch other helpful techniques. Readers are also urged to consult the survey papers [59, 122, 142, 171] and the literature on the EM algorithm for a fuller discussion.

8.4 Allele Frequency Estimation

The ABO and Rh genetic loci are usually typed in matching blood donors to blood recipients. The ABO locus incorporates the three alleles A , B , and O and exhibits the four observable phenotypes A , B , AB , and O . These phenotypes arise because each person inherits two alleles, one from his mother and one from his father, and the alleles A and B are genetically dominant to allele O . Dominance amounts to a masking of the O allele by the presence of an A or B allele. For instance, a person inheriting an A allele from one parent and an O allele from the other parent is said to have genotype A/O and is indistinguishable from a person inheriting an A allele from both parents. This second person has genotype A/A .

The MM algorithm for estimating the population frequencies or proportions of the three alleles involves an interplay between observed phenotypes and underlying unobserved genotypes. As just noted, both genotypes A/O and A/A generate the same phenotype A . Likewise, both genotypes B/O and B/B generate the same phenotype B . Phenotypes AB and O correspond to the single genotypes A/B and O/O , respectively.

As a concrete example, Clarke et al. [50] noted that among their population sample of $n = 521$ duodenal ulcer patients, a total of $n_A = 186$ had phenotype A , $n_B = 38$ had phenotype B , $n_{AB} = 13$ had phenotype AB , and $n_O = 284$ had phenotype O . If we want to estimate the frequencies p_A , p_B , and p_O of the three different alleles from this sample, then we can employ the MM algorithm with the four phenotype counts as the observed data.

The likelihood of the data is given by the multinomial distribution in conjunction with the Hardy-Weinberg law of population genetics. This law specifies that each genotype frequency equals the product of the corresponding allele frequencies with an extra factor of 2 included to account for ambiguity in parental source when the two alleles differ. For example, genotype A/A has frequency p_A^2 , and genotype A/O has frequency $2p_A p_O$.

These assumptions are summarized in the multinomial loglikelihood

$$f(\mathbf{p}) = n_A \ln(p_A^2 + 2p_{AP}p_O) + n_B \ln(p_B^2 + 2p_{BP}p_O) + n_{AB} \ln(2p_{AP}p_B) \\ + n_O \ln p_O^2 + \ln \binom{n}{n_A, n_B, n_{AB}, n_O}.$$

In maximum likelihood estimation we maximize this function of the allele frequencies subject to the equality constraint $p_A + p_B + p_O = 1$ and the nonnegativity constraints $p_A \geq 0$, $p_B \geq 0$, and $p_O \geq 0$.

The loglikelihood function $f(\mathbf{p})$ would be easy to maximize if it were not for the terms $\ln(p_A^2 + 2p_{AP}p_O)$ and $\ln(p_B^2 + 2p_{BP}p_O)$. In the MM algorithm we attack these functions using the convexity of the function $-\ln x$ and the majorization (8.5). This yields the minorization

$$\ln(p_A^2 + 2p_{AP}p_O) \geq \frac{p_{mA}^2}{p_{mA}^2 + 2p_{mAP}p_{mO}} \ln \left(\frac{p_{mA}^2 + 2p_{mAP}p_{mO}}{p_{mA}^2} p_A^2 \right) \\ + \frac{2p_{mAP}p_{mO}}{p_{mA}^2 + 2p_{mAP}p_{mO}} \ln \left(\frac{p_{mA}^2 + 2p_{mAP}p_{mO}}{2p_{mAP}p_{mO}} 2p_{AP}p_O \right).$$

A similar minorization applies to $\ln(p_B^2 + 2p_{BP}p_O)$. These maneuvers have the virtue of separating parameters because logarithms turn products into sums.

Notationally, things become clearer if we introduce the abbreviations

$$n_{mA/A} = n_A \frac{p_{mA}^2}{p_{mA}^2 + 2p_{mAP}p_{mO}} \\ n_{mA/O} = n_A \frac{2p_{mAP}p_{mO}}{p_{mA}^2 + 2p_{mAP}p_{mO}}$$

and likewise for $n_{mB/B}$ and $n_{mB/O}$. We are now faced with maximizing the surrogate function

$$g(\mathbf{p} \mid \mathbf{p}_m) = n_{mA/A} \ln p_A^2 + n_{mA/O} \ln(2p_{AP}p_O) + n_{mB/B} \ln p_B^2 \\ + n_{mB/O} \ln(2p_{BP}p_O) + n_{AB} \ln(2p_{AP}p_B) + n_O \ln p_O^2 + c,$$

where c is an irrelevant constant that depends on the current iterate \mathbf{p}_m but not on the potential value \mathbf{p} of the next iterate. This completes the minorization step of the algorithm.

The maximization step can be accomplished by introducing a Lagrange multiplier and finding a stationary point of the Lagrangian

$$\mathcal{L}(\mathbf{p}, \lambda) = g(\mathbf{p} \mid \mathbf{p}_m) + \lambda(p_A + p_B + p_O - 1).$$

Here we ignore the nonnegativity constraints under the assumption that they are inactive at the solution. Setting the partial derivatives of $\mathcal{L}(\mathbf{p}, \lambda)$,

$$\frac{\partial}{\partial p_A} \mathcal{L}(\mathbf{p}, \lambda) = \frac{2n_{mA/A}}{p_A} + \frac{n_{mA/O}}{p_A} + \frac{n_{AB}}{p_A} + \lambda$$

TABLE 8.1. Iterations for ABO duodenal ulcer data

Iteration m	p_{mA}	p_{mB}	p_{mO}
0	0.3000	0.2000	0.5000
1	0.2321	0.0550	0.7129
2	0.2160	0.0503	0.7337
3	0.2139	0.0502	0.7359
4	0.2136	0.0501	0.7363
5	0.2136	0.0501	0.7363

$$\begin{aligned} \frac{\partial}{\partial p_B} \mathcal{L}(\mathbf{p}, \lambda) &= \frac{2n_{mB/B}}{p_B} + \frac{n_{mB/O}}{p_B} + \frac{n_{AB}}{p_B} + \lambda \\ \frac{\partial}{\partial p_O} \mathcal{L}(\mathbf{p}, \lambda) &= \frac{n_{mA/O}}{p_O} + \frac{n_{mB/O}}{p_O} + \frac{2n_O}{p_O} + \lambda \\ \frac{\partial}{\partial \lambda} \mathcal{L}(\mathbf{p}, \lambda) &= p_A + p_B + p_O - 1, \end{aligned}$$

equal to 0 provides the unique stationary point of $\mathcal{L}(\mathbf{p}, \lambda)$. The solution of the resulting equations is

$$\begin{aligned} p_{m+1,A} &= \frac{2n_{mA/A} + n_{mA/O} + n_{AB}}{2n} \\ p_{m+1,B} &= \frac{2n_{mB/B} + n_{mB/O} + n_{AB}}{2n} \\ p_{m+1,O} &= \frac{n_{mA/O} + n_{mB/O} + 2n_O}{2n}. \end{aligned}$$

In other words, the MM update is identical to a form of gene counting in which the unknown genotype counts are imputed based on the current allele frequency estimates [239]. In these updates, the denominator $2n$ is the total number of genes; the numerators are the current best guesses of the number of alleles of each type contained in the hidden and manifest genotypes.

Table 8.1 shows the progress of the MM iterates starting from the initial estimates $p_{0A} = 0.3$, $p_{0B} = 0.2$, and $p_{0O} = 0.5$. The MM updates are simple enough to carry out on a pocket calculator. Convergence occurs quickly in this example.

8.5 Linear Regression

Because t^2 is a convex function, we can majorize each summand of the sum of squares criterion $\sum_{i=1}^n (y_i - \mathbf{x}_i^* \boldsymbol{\theta})^2$ using inequality (8.6). The overall

surrogate

$$g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = \sum_{i=1}^n \sum_j \alpha_{ij} \left[y_i - \frac{x_{ij}}{\alpha_{ij}} (\theta_j - \theta_{mj}) - \mathbf{x}_i^* \boldsymbol{\theta}_m \right]^2,$$

achieves equality with the sum of squares when $\boldsymbol{\theta} = \boldsymbol{\theta}_m$. Minimization of $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ yields the updates

$$\theta_{m+1,j} = \theta_{mj} + \frac{\sum_{i=1}^n x_{ij} (y_i - \mathbf{x}_i^* \boldsymbol{\theta}_m)}{\sum_{i=1}^n \frac{x_{ij}^2}{\alpha_{ij}}} \quad (8.10)$$

and avoids matrix inversion [171]. Although it seems intuitively reasonable to take $p = 1$ in choosing

$$\alpha_{ij} = \frac{|x_{ij}|^p}{(\sum_k |x_{ik}|^p)},$$

conceivably other values of p might perform better. In fact, it might accelerate convergence to alternate different values of p as the iterations proceed. For problems involving just a few parameters, this iterative scheme is clearly inferior to the usual single-step solution via matrix inversion. Cyclic coordinate descent also avoids matrix operations, and Problem 11 suggests that it will converge faster than the MM update (8.10).

Least squares estimation suffers from the fact that it is strongly influenced by observations far removed from their predicted values. In least absolute deviation regression, we replace $\sum_{i=1}^n (y_i - \mathbf{x}_i^* \boldsymbol{\theta})^2$ by

$$h(\boldsymbol{\theta}) = \sum_{i=1}^n |y_i - \mathbf{x}_i^* \boldsymbol{\theta}| = \sum_{i=1}^n |r_i(\boldsymbol{\theta})|, \quad (8.11)$$

where $r_i(\boldsymbol{\theta}) = y_i - \mathbf{x}_i^* \boldsymbol{\theta}$ is the i th residual. We are now faced with minimizing a nondifferentiable function. Fortunately, the MM algorithm can be implemented by exploiting the concavity of the function \sqrt{u} in inequality (8.7). Because

$$\sqrt{u} \leq \sqrt{u_m} + \frac{u - u_m}{2\sqrt{u_m}}, \quad (8.12)$$

we find that

$$\begin{aligned} h(\boldsymbol{\theta}) &= \sum_{i=1}^n \sqrt{r_i^2(\boldsymbol{\theta})} \\ &\leq h(\boldsymbol{\theta}_m) + \frac{1}{2} \sum_{i=1}^n \frac{r_i^2(\boldsymbol{\theta}) - r_i^2(\boldsymbol{\theta}_m)}{\sqrt{r_i^2(\boldsymbol{\theta}_m)}} \\ &= g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m). \end{aligned}$$

Minimizing $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ is accomplished by minimizing the weighted sum of squares

$$\sum_{i=1}^n w_i(\boldsymbol{\theta}_m) r_i(\boldsymbol{\theta})^2$$

with i th weight $w_i(\boldsymbol{\theta}_m) = |r_i(\boldsymbol{\theta}_m)|^{-1}$. A slight variation of the usual argument for minimizing a sum of squares leads to the update

$$\boldsymbol{\theta}_{m+1} = [\mathbf{X}^* \mathbf{W}(\boldsymbol{\theta}_m) \mathbf{X}]^{-1} \mathbf{X}^* \mathbf{W}(\boldsymbol{\theta}_m) \mathbf{y},$$

where $\mathbf{W}(\boldsymbol{\theta}_m)$ is the diagonal matrix with i th diagonal entry $w_i(\boldsymbol{\theta}_m)$. Unfortunately, the possibility that some weight $w_i(\boldsymbol{\theta}_m)$ is infinite cannot be ruled out. Problem 14 suggests a simple remedy.

8.6 Bradley-Terry Model of Ranking

In the sports version of the Bradley and Terry model [23, 140, 150], each team i in a league of teams is assigned a rank parameter $r_i > 0$. Assuming ties are impossible, team i beats team j with probability $r_i/(r_i + r_j)$. If this outcome occurs y_{ij} times during a season of play, then the probability of the whole season is

$$L(\mathbf{r}) = \prod_{i,j} \left(\frac{r_i}{r_i + r_j} \right)^{y_{ij}},$$

assuming the games are independent. To rank the teams, we find the values \hat{r}_i that maximize $f(\mathbf{r}) = \ln L(\mathbf{r})$. The team with largest \hat{r}_i is considered best, the team with smallest \hat{r}_i is considered worst, and so forth. In view of the fact that $\ln u$ is concave, inequality (8.7) implies

$$\begin{aligned} f(\mathbf{r}) &= \sum_{i,j} y_{ij} \left[\ln r_i - \ln(r_i + r_j) \right] \\ &\geq \sum_{i,j} y_{ij} \left[\ln r_i - \ln(r_{mi} + r_{mj}) - \frac{r_i + r_j - r_{mi} - r_{mj}}{r_{mi} + r_{mj}} \right] \\ &= g(\mathbf{r} \mid \mathbf{r}_m) \end{aligned}$$

with equality when $\mathbf{r} = \mathbf{r}_m$. Differentiating $g(\mathbf{r} \mid \mathbf{r}_m)$ with respect to the i th component r_i of \mathbf{r} and setting the result equal to 0 produces the next iterate

$$r_{m+1,i} = \frac{\sum_{j \neq i} y_{ij}}{\sum_{j \neq i} (y_{ij} + y_{ji}) / (r_{mi} + r_{mj})}.$$

Because $L(\mathbf{r}) = L(\beta \mathbf{r})$ for any $\beta > 0$, we constrain $r_1 = 1$ and omit the update $r_{m+1,1}$. In this example, the MM algorithm separates parameters and allows us to maximize $g(\mathbf{r} \mid \mathbf{r}_m)$ parameter by parameter.

8.7 Linear Logistic Regression

In linear logistic regression, we observe a sequence of independent Bernoulli trials, each resulting in success or failure. The success probability of the i th trial

$$\pi_i(\boldsymbol{\theta}) = \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}}}$$

depends on a predictor (covariate) vector \mathbf{x}_i and parameter vector $\boldsymbol{\theta}$ by analogy with linear regression. The observation y_i at trial i equals 1 for a success and 0 for a failure. In this notation, the likelihood of the data is

$$L(\boldsymbol{\theta}) = \prod_i \pi_i(\boldsymbol{\theta})^{y_i} [1 - \pi_i(\boldsymbol{\theta})]^{1-y_i}.$$

As usual in maximum likelihood estimation, we pass to the loglikelihood

$$f(\boldsymbol{\theta}) = \sum_i [y_i \ln \pi_i(\boldsymbol{\theta}) + (1 - y_i) \ln [1 - \pi_i(\boldsymbol{\theta})].$$

Straightforward calculations show

$$\begin{aligned} df(\boldsymbol{\theta}) &= \sum_i [y_i - \pi_i(\boldsymbol{\theta})] \mathbf{x}_i^* \\ d^2 f(\boldsymbol{\theta}) &= - \sum_i \pi_i(\boldsymbol{\theta}) [1 - \pi_i(\boldsymbol{\theta})] \mathbf{x}_i \mathbf{x}_i^*. \end{aligned}$$

The loglikelihood $f(\boldsymbol{\theta})$ is therefore concave, and we seek to minorize it by a quadratic rather than majorize it by a quadratic as suggested in inequality (8.8). Hence, we must identify a matrix \mathbf{B} such that \mathbf{B} is negative definite and $\mathbf{B} - d^2 f(\boldsymbol{\theta})$ is negative semidefinite for all $\boldsymbol{\theta}$. In view of the scalar inequality $\pi(1 - \pi) \leq \frac{1}{4}$, we take $\mathbf{B} = -\frac{1}{4} \sum_i \mathbf{x}_i \mathbf{x}_i^*$. Maximization of the minorizing quadratic

$$f(\boldsymbol{\theta}_m) + df(\boldsymbol{\theta}_m)(\boldsymbol{\theta} - \boldsymbol{\theta}_m) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_m)^* \mathbf{B}(\boldsymbol{\theta} - \boldsymbol{\theta}_m)$$

is a problem we have met before. It does involve inversion of the matrix \mathbf{B} , but once we have computed \mathbf{B}^{-1} , we can store and reuse it at every iteration.

8.8 Geometric and Signomial Programs

The idea behind these minimization algorithms is best understood in a concrete setting. Consider the posynomial

$$f(\mathbf{x}) = \frac{1}{x_1^3} + \frac{3}{x_1 x_2^2} + x_1 x_2 \tag{8.13}$$

with the implied constraints $x_1 > 0$ and $x_2 > 0$. The majorization (8.9) applied to the third term of $f(\mathbf{x})$ yields

$$\begin{aligned} x_1 x_2 &\leq x_{m1} x_{m2} \left[\frac{1}{2} \left(\frac{x_1}{x_{m1}} \right)^2 + \frac{1}{2} \left(\frac{x_2}{x_{m2}} \right)^2 \right] \\ &= \frac{x_{m2}}{2x_{m1}} x_1^2 + \frac{x_{m1}}{2x_{m2}} x_2^2. \end{aligned}$$

Applied to the second term of $f(\mathbf{x})$, it gives with x_1^{-1} replacing x_1 and x_2^{-1} replacing x_2 ,

$$\begin{aligned} \frac{3}{x_1 x_2^2} &\leq \frac{3}{x_{m1} x_{m2}^2} \left[\frac{1}{3} \left(\frac{x_{m1}}{x_1} \right)^3 + \frac{2}{3} \left(\frac{x_{m2}}{x_2} \right)^3 \right] \\ &= \frac{x_{m1}^2}{x_{m2}^2} \frac{1}{x_1^3} + \frac{2x_{m2}}{x_{m1}} \frac{1}{x_2^3}. \end{aligned}$$

The second step of the MM algorithm for minimizing $f(\mathbf{x})$ therefore splits into minimizing the two surrogate functions

$$\begin{aligned} g_1(x_1 \mid \mathbf{x}_m) &= \frac{1}{x_1^3} + \frac{x_{m1}^2}{x_{m2}^2} \frac{1}{x_1^3} + \frac{x_{m2}}{2x_{m1}} x_1^2 \\ g_2(x_2 \mid \mathbf{x}_m) &= \frac{2x_{m2}}{x_{m1}} \frac{1}{x_2^3} + \frac{x_{m1}}{2x_{m2}} x_2^2. \end{aligned}$$

If we set the derivatives of each of these equal to 0, then we find the solutions

$$\begin{aligned} x_{m+1,1} &= \sqrt[5]{3 \left(\frac{x_{m1}^2}{x_{m2}^2} + 1 \right) \frac{x_{m1}}{x_{m2}}} \\ x_{m+1,2} &= \sqrt[5]{6 \frac{x_{m2}^2}{x_{m1}}}. \end{aligned}$$

It is obvious that the point $\mathbf{x} = (\sqrt[5]{6}, \sqrt[5]{6})^*$ is a fixed point of these equations and minimizes $f(\mathbf{x})$. Ignoring this fact, Table 8.2 records the iterates of both the MM algorithm and cyclic coordinate descent. Although the MM updates are slower to converge, they are less complicated. See Problem 4 of Chap. 7 for the form of the cyclic coordinate descent updates.

This MM analysis carries over to general posynomials except that we cannot expect to derive explicit solutions of the minimization step. (See Problem 28.) Each separated surrogate function is a posynomial in a single variable. If the powers appearing in one of these posynomials are integers, then the derivative of the posynomial is a rational function, and once we equate it to 0, we are faced with solving a polynomial equation. This can be accomplished by bisection or by Newton's method as discussed in Chap. 10. Introducing posynomial constraints is another matter. Box constraints in

TABLE 8.2. MM and coordinate descent iterates for a geometric program

m	MM algorithm			Coordinate descent		
	x_{m1}	x_{m2}	$f(\mathbf{x}_m)$	x_{m1}	x_{m2}	$f(\mathbf{x}_m)$
0	1.00000	2.00000	3.75000	1.00000	2.00000	3.75000
1	1.13397	1.88818	3.56899	1.19437	1.61420	3.47886
2	1.19643	1.75472	3.49766	1.32882	1.50339	3.42280
3	1.24544	1.66786	3.46079	1.38616	1.46165	3.41457
4	1.28395	1.60829	3.44074	1.41117	1.44432	3.41312
5	1.31428	1.56587	3.42942	1.42219	1.43685	3.41285
10	1.39358	1.47003	3.41427	1.43082	1.43107	3.41279
20	1.42699	1.43496	3.41280	1.43097	1.43097	3.41279
30	1.43054	1.43140	3.41279	1.43097	1.43097	3.41279
40	1.43092	1.43101	3.41279	1.43097	1.43097	3.41279
50	1.43096	1.43097	3.41279	1.43097	1.43097	3.41279
51	1.43097	1.43097	3.41279	1.43097	1.43097	3.41279

the form $a_i \leq x_i \leq b_i$ are consistent with parameter separation as developed here, but more complicated posynomial constraints are not.

The perturbation

$$f(\mathbf{x}) = \frac{1}{x_1^3} + \frac{3}{x_1 x_2^2} + x_1 x_2 - \sqrt{x_1 x_2}$$

of the function (8.13) is called a signomial rather than a posynomial [20]. If we want to minimize this new function subject to the constraints $x_1 > 0$ and $x_2 > 0$, then a different tactic is needed for the terms with negative coefficients [170]. Consider the minorization $z \geq 1 + \ln z$ around the point $z = 1$ derived from the convexity of $-\ln z$. If we let $z = \sqrt{x_1 x_2} / \sqrt{x_{m1} x_{m2}}$, then the more elaborate minorization

$$\sqrt{x_1 x_2} \geq \frac{1}{2} \sqrt{x_{m1} x_{m2}} (2 + \ln x_1 + \ln x_2 - \ln x_{m1} - \ln x_{m2})$$

follows. Multiplication of this by -1 now gives the operative majorization. Up to an irrelevant additive constant, the overall surrogate function equals the sum of the two parameter separated surrogates

$$g_1(x_1 | \mathbf{x}_m) = \frac{1}{x_1^3} + \frac{x_{m1}^2}{x_m^2} \frac{1}{x_1^3} + \frac{x_{m2}}{2x_{m1}} x_1^2 - \frac{1}{2} \sqrt{x_{m1} x_{m2}} \ln x_1$$

$$g_2(x_2 | \mathbf{x}_m) = \frac{2x_{m2}}{x_{m1}} \frac{1}{x_2^3} + \frac{x_{m1}}{2x_{m2}} x_2^2 - \frac{1}{2} \sqrt{x_{m1} x_{m2}} \ln x_2.$$

The corresponding minima are roots of the polynomial equations

$$0 = \frac{x_{m2}}{x_{m1}} x_1^5 - \frac{1}{2} \sqrt{x_{m1} x_{m2}} x_1^3 - 3 \left(1 + \frac{x_{m1}^2}{x_{m2}^2} \right)$$

$$0 = \frac{x_{m1}}{x_{m2}} x_2^5 - \frac{1}{2} \sqrt{x_{m1} x_{m2}} x_2^3 - \frac{6x_{m2}}{x_{m1}}.$$

Table 8.3 displays the MM iterates for this signomial program. Each update relies on Newton's method to solve the two preceding quintic equations. Although there is no convexity guarantee that the converged point is optimal, random sampling of the objective function does not produce a better point.

TABLE 8.3. MM iterates for a signomial program

m	x_{m1}	x_{m2}	$f(\mathbf{x}_m)$
0	1.00000	2.00000	2.33579
1	1.19973	2.04970	2.06522
5	1.39153	1.73281	1.94756
10	1.48648	1.61186	1.92935
15	1.52318	1.57173	1.92702
20	1.53763	1.55678	1.92668
25	1.54336	1.55097	1.92663
30	1.54565	1.54867	1.92662
35	1.54656	1.54776	1.92662
40	1.54692	1.54740	1.92662
45	1.54706	1.54725	1.92662
50	1.54712	1.54720	1.92662
55	1.54714	1.54717	1.92662
60	1.54715	1.54716	1.92662
65	1.54716	1.54716	1.92662

8.9 Poisson Processes

In preparation for our exposition of transmission tomography in the next section, let us briefly review the theory of Poisson processes, a topic from probability of considerable interest in its own right. A Poisson process involves points randomly scattered in a region S of \mathbb{R}^n [113, 133, 148, 154]. The notion that the points are concentrated on average more in some regions than in others is captured by postulating an intensity function $\lambda(\mathbf{x}) \geq 0$ on S . The expected number of points in a subregion T is given by the integral $\omega = \int_T \lambda(\mathbf{x}) d\mathbf{x}$. If $\omega = \infty$, then an infinite number of random points occur in T . If $\omega < \infty$, then a finite number of random points occur in T , and the probability that this number equals k is given by the Poisson probability

$$p_k(\omega) = \frac{\omega^k}{k!} e^{-\omega}.$$

Derivation of this formula depends critically on the assumption that the numbers N_{T_i} of random points in disjoint regions T_i are independent

random variables. This basically means that knowing the values of some of the N_{T_i} tells one nothing about the values of the remaining N_{T_i} . The model also presupposes that random points never coincide.

The Poisson distribution has a peculiar relationship to the multinomial distribution. Suppose a Poisson random variable Z with mean ω represents the number of outcomes from some experiment, say an experiment involving a Poisson process. Let each outcome be independently classified in one of l categories, the k th of which occurs with probability p_k . Then the number of outcomes Z_k falling in category k is Poisson distributed with mean $\omega_k = p_k\omega$. Furthermore, the random variables Z_1, \dots, Z_l are independent. Conversely, if $Z = \sum_{k=1}^l Z_k$ is a sum of independent Poisson random variables Z_k with means $\omega_k = p_k\omega$, then conditional on $Z = n$, the vector $(Z_1, \dots, Z_l)^*$ follows a multinomial distribution with n trials and cell probabilities p_1, \dots, p_l . To prove the first two of these assertions, let $n = n_1 + \dots + n_l$. Then

$$\begin{aligned} \Pr(Z_1 = n_1, \dots, Z_l = n_l) &= \frac{\omega^n}{n!} e^{-\omega} \binom{n}{n_1, \dots, n_l} \prod_{k=1}^l p_k^{n_k} \\ &= \prod_{k=1}^l \frac{\omega_k^{n_k}}{n_k!} e^{-\omega_k} \\ &= \prod_{k=1}^l \Pr(Z_k = n_k). \end{aligned}$$

To prove the converse, divide the last string of equalities by the probability $\Pr(Z = n) = \omega^n e^{-\omega} / n!$.

The random process of assigning points to categories is termed coloring in the stochastic process literature. When there are just two colors, and only random points of one of the colors are tracked, then the process is termed random thinning. We will see examples of both coloring and thinning in the next section.

8.10 Transmission Tomography

Problems in medical imaging often involve thousands of parameters. As an illustration of the MM algorithm, we treat maximum likelihood estimation in transmission tomography. Traditionally, transmission tomography images have been reconstructed by the methods of Fourier analysis. Fourier methods are fast but do not take into account the uncertainties of photon counts. Statistically based methods give better reconstructions with less patient exposure to harmful radiation.

The purpose of transmission tomography is to reconstruct the local attenuation properties of the object being imaged [124]. Attenuation is to

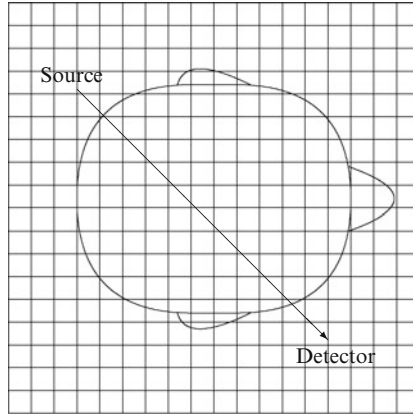


FIGURE 8.2. Cartoon of transmission tomography

be roughly equated with density. In medical applications, material such as bone is dense and stops or deflects X-rays (high-energy photons) better than soft tissue. With enough photons, even small gradations in soft tissue can be detected. A two-dimensional image is constructed from a sequence of photon counts. Each count corresponds to a projection line L drawn from an X-ray source through the imaged object to an X-ray detector. The average number of photons sent from the source along L to the detector is known in advance. The random number of photons actually detected is determined by the probability of a single photon escaping deflection or capture along L . Figure 8.2 shows one such projection line beamed through a cartoon of the human head.

To calculate this probability, we let $\mu(\mathbf{x})$ be the intensity (or attenuation coefficient) of photon deflection or capture per unit length at the point $\mathbf{x} = (x_1, x_2)$ in the plane. We can imagine that deflection or capture events occur completely randomly along L according to a Poisson process. The first such event effectively prevents the photon from being detected. Thus, the photon is detected with the Poisson probability $p_0(\omega) = e^{-\omega}$ of no such events, where

$$\omega = \int_L \mu(\mathbf{x}) ds$$

is the line integral of $\mu(\mathbf{x})$ along L . In actual practice, X-rays are beamed through the object along a large number of different projection lines. We therefore face the inverse problem of reconstructing a function $\mu(\mathbf{x})$ in the plane from a large number of its measured line integrals. Imposing enough smoothness on $\mu(\mathbf{x})$, one can solve this classical deterministic problem by applying Radon transform techniques from Fourier analysis [124].

An alternative to the Fourier method is to pose an explicitly stochastic model and estimate its parameters by maximum likelihood [167, 168].

The MM algorithm suggests itself in this context. The stochastic model depends on dividing the object of interest into small nonoverlapping regions of constant attenuation called pixels. Typically the pixels are squares on a regular grid as depicted in Fig. 8.2. The attenuation attributed to pixel j constitutes parameter θ_j of the model. Since there may be thousands of pixels, implementation of maximum likelihood algorithms such as scoring or Newton's method as discussed in Chap. 10 is out of the question.

To summarize our discussion, each observation Y_i is generated by beaming a stream of X-rays or high-energy photons from an X-ray source toward some detector on the opposite side of the object. The observation (or projection) Y_i counts the number of photons detected along the i th line of flight. Naturally, only a fraction of the photons are successfully transmitted from source to detector. If l_{ij} is the length of the segment of projection line i intersecting pixel j , then the probability of a photon escaping attenuation along projection line i is the exponentiated line integral $\exp(-\sum_j l_{ij}\theta_j)$.

In the absence of the intervening object, the number of photons generated and ultimately detected follows a Poisson distribution. We assume that the mean d_i of this distribution for projection line i is known. Ideally, detectors are long tubes aimed at the source. If a photon is deflected, then it is detected neither by the tube toward which it is initially headed nor by any other tube. In practice, many different detectors collect photons simultaneously from a single source. If we imagine coloring the tubes, then each photon is colored by the tube toward which it is directed. Each stream of colored photons is then thinned by capture or deflection. These considerations imply that the counts Y_i are independent and Poisson distributed with means $d_i \exp(-\sum_j l_{ij}\theta_j)$. It follows that we can express the loglikelihood of the observed data $Y_i = y_i$ as the finite sum

$$\sum_i \left[-d_i e^{-\sum_j l_{ij}\theta_j} - y_i \sum_j l_{ij}\theta_j + y_i \ln d_i - \ln y_i! \right]. \quad (8.14)$$

Omitting irrelevant constants, we can rewrite the loglikelihood (8.14) more succinctly as

$$L(\boldsymbol{\theta}) = -\sum_i f_i(\mathbf{l}_i^* \boldsymbol{\theta}),$$

where $f_i(t)$ is the convex function $d_i e^{-t} + y_i t$ and $\mathbf{l}_i^* \boldsymbol{\theta} = \sum_j l_{ij}\theta_j$ is the inner product of the attenuation parameter vector $\boldsymbol{\theta}$ and the vector of intersection lengths \mathbf{l}_i for projection i .

To generate a surrogate function, we majorize each $f_i(\mathbf{l}_i^* \boldsymbol{\theta})$ according to the recipe (8.5). This gives the surrogate function

$$g(\boldsymbol{\theta} | \boldsymbol{\theta}_m) = -\sum_i \sum_j \frac{l_{ij}\theta_{m,j}}{\mathbf{l}_i^* \boldsymbol{\theta}_m} f_i\left(\frac{\mathbf{l}_i^* \boldsymbol{\theta}_m}{\theta_{m,j}} \theta_j\right) \quad (8.15)$$

minorizing $L(\boldsymbol{\theta})$. By construction, maximization of $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ separates into a sequence of one-dimensional problems, each of which can be solved approximately by one step of Newton's method. We will take up the details of this in Chap. 10.

The images produced by maximum likelihood estimation in transmission tomography look grainy. The cure is to enforce image smoothness by penalizing large differences between estimated attenuation parameters of neighboring pixels. Geman and McClure [103] recommend multiplying the likelihood of the data by a Gibbs prior $\pi(\boldsymbol{\theta})$. Equivalently we add the log prior

$$\ln \pi(\boldsymbol{\theta}) = -\gamma \sum_{\{j,k\} \in N} w_{jk} \psi(\theta_j - \theta_k)$$

to the loglikelihood, where γ and the weights w_{jk} are positive constants, N is a set of unordered pairs $\{j, k\}$ defining a neighborhood system, and $\psi(r)$ is called a potential function. This function should be large whenever $|r|$ is large. Neighborhoods have limited extent. For instance, if the pixels are squares, we might define the weights by $w_{jk} = 1$ for orthogonal nearest neighbors sharing a side and $w_{jk} = 1/\sqrt{2}$ for diagonal nearest neighbors sharing only a corner. The constant γ scales the overall strength assigned to the prior. The sum $L(\boldsymbol{\theta}) + \ln \pi(\boldsymbol{\theta})$ is called the log posterior function; its maximum is the posterior mode.

Choice of the potential function $\psi(r)$ is the most crucial feature of the Gibbs prior. It is convenient to assume that $\psi(r)$ is even and strictly convex. Strict convexity leads to the strict concavity of the log posterior function $L(\boldsymbol{\theta}) + \ln \pi(\boldsymbol{\theta})$ and permits simple modification of the MM algorithm based on the surrogate function $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ defined by equation (8.15). Many potential functions exist satisfying these conditions. One natural example is $\psi(r) = r^2$. This choice unfortunately tends to deter the formation of boundaries. The gentler alternatives $\psi(r) = \sqrt{r^2 + \epsilon}$ for a small positive ϵ and $\psi(r) = \ln[\cosh(r)]$ are preferred in practice [111]. Problem 35 asks the reader to verify some of the properties of these two potential functions.

One adverse consequence of introducing a prior is that it couples pairs of parameters in the maximization step of the MM algorithm for finding the posterior mode. One can decouple the parameters by exploiting the convexity and evenness of the potential function $\psi(r)$ through the inequality

$$\begin{aligned} \psi(\theta_j - \theta_k) &= \psi\left(\frac{1}{2}[2\theta_j - \theta_{mj} - \theta_{mk}] + \frac{1}{2}[-2\theta_k + \theta_{mj} + \theta_{mk}]\right) \\ &\leq \frac{1}{2}\psi(2\theta_j - \theta_{mj} - \theta_{mk}) + \frac{1}{2}\psi(2\theta_k - \theta_{mj} - \theta_{mk}), \end{aligned}$$

which is strict unless $\theta_j + \theta_k = \theta_{mj} + \theta_{mk}$. This inequality allows us to redefine the surrogate function as

$$g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$$

$$\begin{aligned}
 &= - \sum_i \sum_j \frac{l_{ij} \theta_{mj}}{\mathbf{l}_i^* \boldsymbol{\theta}_m} f_i \left(\frac{\mathbf{l}_i^* \boldsymbol{\theta}_m}{\theta_{mj}} \theta_j \right) \\
 &\quad - \frac{\gamma}{2} \sum_{\{j,k\} \in N} w_{jk} [\psi(2\theta_j - \theta_{mj} - \theta_{mk}) + \psi(2\theta_k - \theta_{mj} - \theta_{mk})].
 \end{aligned}$$

Once again the parameters are separated, and the maximization step reduces to a sequence of one-dimensional problems. Maximizing $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ drives the log posterior uphill and eventually leads to the posterior mode.

8.11 Poisson Multigraphs

In a graph the number of edges between any two nodes is 0 or 1. A multigraph allows an arbitrary number of edges between any two nodes. Multigraphs are natural structures for modeling the internet and gene and protein networks. Here we consider a multigraph with a random number of edges X_{ij} connecting every pair of nodes $\{i, j\}$. In particular, we assume that the X_{ij} are independent Poisson random variables with means μ_{ij} . As a plausible model for ranking nodes, we take $\mu_{ij} = p_i p_j$, where p_i and p_j are nonnegative propensities [217].

The loglikelihood of the observed edge counts $x_{ij} = x_{ji}$ amounts to

$$\begin{aligned}
 L(\mathbf{p}) &= \sum_{\{i,j\}} (x_{ij} \ln \mu_{ij} - \mu_{ij} - \ln x_{ij}!) \\
 &= \sum_{\{i,j\}} [x_{ij} (\ln p_i + \ln p_j) - p_i p_j - \ln x_{ij}!].
 \end{aligned}$$

Inspection of $L(\mathbf{p})$ shows that the parameters are separated except for the products $p_i p_j$. To achieve full separation of parameters in maximum likelihood estimation, we employ the majorization

$$p_i p_j \leq \frac{p_{mj}}{2p_{mi}} p_i^2 + \frac{p_{mi}}{2p_{mj}} p_j^2.$$

Equality prevails here when $\mathbf{p} = \mathbf{p}_m$. This majorization leads to the minorization

$$\begin{aligned}
 L(\mathbf{p}) &\geq \sum_{\{i,j\}} [x_{ij} (\ln p_i + \ln p_j) - \frac{p_{mj}}{2p_{mi}} p_i^2 - \frac{p_{mi}}{2p_{mj}} p_j^2 - \ln x_{ij}!] \\
 &= g(\mathbf{p} \mid \mathbf{p}_m).
 \end{aligned}$$

Maximization of $g(\mathbf{p} \mid \mathbf{p}_m)$ can be accomplished by setting

$$\frac{\partial}{\partial p_i} g(\mathbf{p} \mid \mathbf{p}_m) = \sum_{j \neq i} \frac{x_{ij}}{p_i} - \sum_{j \neq i} \frac{p_{mj}}{p_{mi}} p_i = 0$$

The solution

$$p_{m+1,i} = \sqrt{\frac{p_{mi} \sum_{j \neq i} x_{ij}}{\sum_{j \neq i} p_{mj}}} \quad (8.16)$$

is straightforward to implement and maps positive parameters to positive parameters. When edges are sparse, the range of summation in $\sum_{j \neq i} x_{ij}$ can be limited to those nodes j with $x_{ij} > 0$. Observe that these sums need only be computed once. The partial sums $\sum_{j \neq i} p_{mj} = \sum_j p_{mj} - p_{mi}$ require updating the full sum $\sum_j p_{mj}$ once per iteration.

A similar MM algorithm can be derived for a Poisson model of arc formation in a directed multigraph. We now postulate a donor propensity p_i and a recipient propensity q_j for arcs extending from node i to node j . If the number of such arcs X_{ij} is Poisson distributed with mean $p_i q_j$, then under independence we have the loglikelihood

$$L(\mathbf{p}, \mathbf{q}) = \sum_i \sum_{j \neq i} [x_{ij}(\ln p_i + \ln q_j) - p_i q_j - \ln x_{ij}!]$$

With directed arcs the observed numbers x_{ij} and x_{ji} may differ. The minorization

$$L(\mathbf{p}, \mathbf{q}) \geq \sum_i \sum_{j \neq i} [x_{ij}(\ln p_i + \ln q_j) - \frac{q_{mj}}{2p_{mi}} p_i^2 - \frac{p_{mi}}{2q_{mj}} q_j^2 - \ln x_{ij}!]$$

now yields the MM updates

$$p_{m+1,i} = \sqrt{\frac{p_{mi} \sum_{j \neq i} x_{ij}}{\sum_{j \neq i} q_{mj}}}, \quad q_{m+1,j} = \sqrt{\frac{q_{mj} \sum_{i \neq j} x_{ij}}{\sum_{i \neq j} p_{mi}}}. \quad (8.17)$$

Again these are computationally simple to implement and map positive parameters to positive parameters. It is important to observe that the loglikelihood $L(\mathbf{p}, \mathbf{q})$ is invariant under the rescaling cp_i and $c^{-1}q_j$ for some positive constant c and all i and j . This fact suggests that we fix one propensity, say $p_1 = 1$, and omit its update.

There are interesting examples of propensity estimation in literary analysis and attribution. Nodes are words in a text. An arc is drawn between two consecutive words, from the first word to the second word, provided the words are not separated by a punctuation mark. Here we examine Charlotte Bronte's novel *Jane Eyre*. Based on the directed multigraph model, it is possible to calculate the donor and recipient propensities of each word. Given these propensities, one can assign a p-value under the Poisson distribution indicating whether two words have an excess number of connections. Table 8.4 lists the most significant ordered word pairs in *Jane Eyre*.

TABLE 8.4. Most significantly connected word pairs in *Jane Eyre*

Rank	-Log p-value	Observed	Expected	Pair
1	-304.77	323	14.30	I am
2	-293.43	420	33.82	It was
3	-271.98	510	62.60	I had
4	-258.13	306	17.28	It is
5	-256.50	251	9.24	You are
6	-239.92	811	192.31	Of the
7	-233.92	155	1.82	Do not
8	-219.41	609	119.73	In the
9	-208.02	173	4.17	Could not
10	-196.42	320	32.02	To be
11	-191.60	154	3.37	Had been
12	-179.88	259	21.18	I could
13	-168.41	138	3.20	Did not
14	-162.75	337	47.45	On the
15	-153.04	317	44.68	I have
16	-152.62	480	106.93	I was
17	-132.87	108	2.46	Have been
18	-118.87	162	12.10	I should
19	-117.57	112	3.91	As if
20	-115.97	130	6.61	There was
21	-114.57	110	3.92	Would be
22	-113.61	123	5.80	Do you
23	-106.61	171	16.85	You have
24	-104.87	62	0.49	At least
25	-103.23	132	8.78	A little

8.12 Problems

1. Prove that the majorization relation between functions is closed under the formation of sums, nonnegative products, limits, and composition with an increasing function. In what sense is the relation also transitive?
2. Demonstrate the majorizing and minorizing inequalities

$$x^q \leq qx_m^{q-1}x + (1-q)x_m^q$$

$$\ln x \leq \frac{x}{x_m} + \ln x_m - 1$$

$$\begin{aligned}
x \ln x &\leq \frac{x^2}{x_m} + x \ln x_m - x \\
\|\mathbf{x}\| &\geq \frac{\mathbf{x}_m^* \mathbf{x}}{\|\mathbf{x}_m\|} \\
xy &\leq \frac{y_m}{2x_m} x^2 + \frac{x_m}{2y_m} y^2 \\
-xy &\leq -x_m y_m \left[1 + \ln \left(\frac{x}{x_m} \right) + \ln \left(\frac{y}{y_m} \right) \right] \\
\frac{1}{x} &\leq \frac{1}{x_m} - \frac{x - x_m}{x_m^2} + \frac{(x - x_m)^2}{c^3} \\
\frac{1}{x + y} &\leq \left(\frac{x_m}{x_m + y_m} \right)^2 \frac{1}{x} + \left(\frac{y_m}{x_m + y_m} \right)^2 \frac{1}{y}
\end{aligned}$$

Determine the relevant domains of each variable q , x , x_m , y , y_m , and c , and check that equality occurs in each of the inequalities when $x = x_m$ and $y = y_m$ [122].

3. As alternatives to the fifth and sixth examples of Problem 2, demonstrate the majorizations

$$\begin{aligned}
xy &\leq \frac{1}{2}(x^2 + y^2) + \frac{1}{2}(x_m - y_m)^2 - (x_m - y_m)(x - y) \\
-xy &\leq \frac{1}{2}(x^2 + y^2) + \frac{1}{2}(x_m + y_m)^2 - (x_m + y_m)(x + y)
\end{aligned}$$

valid for all values of x , y , x_m , and y_m .

4. Based on Problem 3, devise an MM algorithm to minimize Rosenbrock's function

$$f(\mathbf{x}) = 100(x_1^2 - x_2)^2 + (x_1 - 1)^2.$$

Show that up to an irrelevant constant $f(\mathbf{x})$ is majorized by the sum of the two functions

$$\begin{aligned}
g_1(x_1 | \mathbf{x}_m) &= 200x_1^4 - [200(x_{m1}^2 + x_{m2}) - 1]x_1^2 - 2x_1 \\
g_2(x_2 | \mathbf{x}_m) &= 200x_2^2 - 200(x_{m1}^2 + x_{m2})x_2.
\end{aligned}$$

Hence at each iteration one must minimize a quartic in x_1 and a quadratic in x_2 . Implement this MM algorithm, and check whether it converges to the global minimum of $f(\mathbf{x})$ at $\mathbf{x} = \mathbf{1}$.

5. Devise an MM algorithm to minimize Snell's objective function (1.1).
 6. Prove van Ruitenburg's [263] minorization

$$\ln x \geq -\frac{3x_m}{2x} - \frac{x}{2x_m} + \ln x_m + 2.$$

Deduce the further minorization

$$x \ln x \geq -\frac{3x_m}{2} - \frac{x^2}{2x_m} + x \ln x_m + 2x.$$

7. Suppose $p \in [1, 2]$ and $x_m \neq 0$. Verify the majorizing inequality

$$|x|^p \leq \frac{p}{2}|x_m|^{p-2}x^2 + \left(1 - \frac{p}{2}\right)|x_m|^p.$$

This is important in least ℓ_p regression.

8. The majorization

$$|x| \leq \frac{1}{2|x_m|}(x^2 + x_m^2) \quad (8.18)$$

for $x_m \neq 0$ is a special case of Problem 7. Use the majorization (8.18) and the identity

$$\max\{x, y\} = \frac{1}{2}|x - y| + \frac{1}{2}x + \frac{1}{2}y$$

to majorize $\max\{x, y\}$ when $x_m \neq y_m$. Note that your majorization contains the product xy up to a negative factor. Describe how one can invoke Problem 2 or Problem 3 to separate x and y .

9. Prove the majorization (8.9) of the text.

10. Consider the function

$$f(x) = \frac{1}{4}x^4 - \frac{1}{2}x^2.$$

This function has global minima at $x = \pm 1$ and a local maximum at $x = 0$. Show that the function

$$g(x | x_m) = \frac{1}{4}x^4 + \frac{1}{2}x_m^2 - xx_m$$

majorizes $f(x)$ at x_m and leads to the MM update $x_{m+1} = \sqrt[3]{x_m}$. Prove that the alternative update $x_{m+1} = -\sqrt[3]{x_m}$ leads to the same value of $f(x)$, but the first update always converges while the second oscillates in sign and has two converging subsequences [60].

11. In the regression algorithm (8.10), let p tend to 0. If there are q predictors and all x_{ij} are nonzero, then show that $\alpha_{ij} = 1/q$. This leads to the update

$$\theta_{m+1,j} = \theta_{mj} + \frac{\sum_{i=1}^n x_{ij}(y_i - \mathbf{x}_i^* \boldsymbol{\theta}_m)}{q \sum_{i=1}^n x_{ij}^2}. \quad (8.19)$$

On the other hand, argue that cyclic coordinate descent yields the update

$$\theta_{m+1,j} = \theta_{mj} + \frac{\sum_{i=1}^n x_{ij}(y_i - \mathbf{x}_i^* \boldsymbol{\theta}_m)}{\sum_{i=1}^n x_{ij}^2},$$

which definitely takes larger steps.

12. A number μ is said to be a q quantile of the n numbers x_1, \dots, x_n if it satisfies

$$\frac{1}{n} \sum_{x_i \leq \mu_q} 1 \geq q \quad \text{and} \quad \frac{1}{n} \sum_{x_i \geq \mu_q} 1 \geq 1 - q.$$

If we define

$$\rho_q(r) = \begin{cases} qr & r \geq 0 \\ -(1-q)r & r < 0, \end{cases}$$

then demonstrate that μ is a q quantile if and only if μ minimizes the function $f_q(\mu) = \sum_{i=1}^n \rho_q(x_i - \mu)$. Medians correspond to the case $q = 1/2$.

13. Continuing Problem 12, show that the function $\rho_q(r)$ is majorized by the quadratic

$$\zeta_q(r | r_m) = \frac{1}{4} \left[\frac{r^2}{|r_m|} + (4q - 2)r + |r_m| \right].$$

Deduce from this majorization the MM algorithm

$$\begin{aligned} \mu_{m+1} &= \frac{n(2q - 1) + \sum_{i=1}^n w_{mi} x_i}{\sum_{i=1}^n w_{mi}} \\ w_{mi} &= \frac{1}{|x_i - \mu_m|} \end{aligned}$$

for finding a q quantile. This interesting algorithm involves no sorting, only arithmetic operations.

14. Suppose we minimize the function

$$h_\epsilon(\boldsymbol{\theta}) = \sum_{i=1}^p \left\{ \left[y_i - \sum_{j=1}^q x_{ij} \theta_j \right]^2 + \epsilon \right\}^{1/2} \quad (8.20)$$

instead of the function $h(\boldsymbol{\theta})$ in equation (8.11) for a small, positive number ϵ . Show that the same MM algorithm applies with revised weights $w_i(\boldsymbol{\theta}_m) = 1/\sqrt{r_i(\boldsymbol{\theta}_m)^2 + \epsilon}$.

15. At the point \mathbf{y} suppose the affine function $\mathbf{v}_j^* \mathbf{x} + a_j$ is the only contribution to the convex function

$$f(\mathbf{x}) = \max_{1 \leq i \leq n} (\mathbf{v}_i^* \mathbf{x} + a_i).$$

achieving the maximum. Show that the quadratic function

$$\begin{aligned} g(\mathbf{x} \mid \mathbf{y}) &= \frac{c}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \mathbf{v}_j^* \mathbf{x} + a_j \\ &= \frac{c}{2} \left\| \mathbf{x} - \mathbf{y} + \frac{1}{c} \mathbf{v}_j \right\|^2 + \mathbf{v}_j^* \mathbf{y} - \frac{1}{2c} \|\mathbf{v}_j\|^2 + a_j \end{aligned}$$

majorizes $f(\mathbf{x})$ with anchor \mathbf{y} for $c > 0$ sufficiently large. Argue that the best choice of $g(\mathbf{x} \mid \mathbf{y})$ takes $c = \max_{i \neq j} c_i$, where

$$c_i = \frac{\|\mathbf{v}_i - \mathbf{v}_j\|^2}{2[a_j - a_i + (\mathbf{v}_j - \mathbf{v}_i)^* \mathbf{y}]},$$

provided none of the denominators vanish. (Hint: Investigate the tangency conditions between $g(\mathbf{x} \mid \mathbf{y})$ and $\mathbf{v}_i^* \mathbf{x} + a_i$.)

16. Based on Problem 15, consider majorization of the ℓ_∞ norm $\|\mathbf{x}\|_\infty$ around $\mathbf{y} \neq \mathbf{0}$. If j is the sole index with $|y_j| = \|\mathbf{y}\|_\infty$, then prove that

$$\begin{aligned} g(\mathbf{x} \mid \mathbf{y}) &= \frac{c}{2} \|\mathbf{x} - \mathbf{y}\|^2 + \operatorname{sgn}(y_j) x_j \\ c &= \frac{1}{|y_j| - \max_{i \neq j} |y_i|} \end{aligned}$$

majorizes $\|\mathbf{x}\|_\infty$ around \mathbf{y} .

17. Suppose the differentiable function $f(\mathbf{x})$ satisfies

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq L \|\mathbf{y} - \mathbf{x}\|$$

for all \mathbf{y} and \mathbf{x} . Deduce the majorization

$$f(\mathbf{x}) \leq f(\mathbf{x}_m) + df(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) + \frac{L}{2} \|\mathbf{x} - \mathbf{x}_m\|^2$$

from the expansion $f(\mathbf{x}) = f(\mathbf{x}_m) + \int_0^1 df[\mathbf{x}_m + t(\mathbf{x} - \mathbf{x}_m)](\mathbf{x} - \mathbf{x}_m) dt$. Verify the special case

$$\begin{aligned} \|M\mathbf{x} - \mathbf{y}\|^2 &\leq \|M\mathbf{x}_m - \mathbf{y}\|^2 + \|M\|^2 [\|\mathbf{x} - \mathbf{x}_m + \mathbf{z}_m\|^2 - \|\mathbf{z}_m\|^2] \\ \mathbf{z}_m &= \|M\|^{-2} M^*(M\mathbf{x}_m - \mathbf{y}) \end{aligned}$$

for the Euclidean and spectral norms.

18. Validate the Euclidean norm majorization

$$\|\mathbf{y} - \mathbf{x}\|_{\dagger} \leq \|\mathbf{y} - \mathbf{x}_m\|_{\dagger} + c\|\mathbf{x}_m - \mathbf{x}\|.$$

Here $\|\cdot\|_{\dagger}$ an arbitrary norm on \mathbb{R}^n , and c is the positive constant guaranteed by Example 2.5.6.

19. Let \mathbf{P} be an orthogonal projection onto a subspace of \mathbb{R}^n . Demonstrate the majorization

$$\|\mathbf{P}\mathbf{x}\|^2 \leq \|\mathbf{x} - \mathbf{x}_m\|^2 + 2(\mathbf{x} - \mathbf{x}_m)^* \mathbf{P}\mathbf{x}_m + \|\mathbf{P}\mathbf{x}_m\|^2.$$

20. Consider minimizing the function

$$f(\boldsymbol{\mu}) = \sum_{i=1}^p \|\mathbf{x}_i - \boldsymbol{\mu}\|$$

for p points $\mathbf{x}_1, \dots, \mathbf{x}_p$ in \mathbb{R}^n . Demonstrate that

$$g(\boldsymbol{\mu} \mid \boldsymbol{\mu}_m) = \frac{1}{2} \sum_{i=1}^p \frac{\|\mathbf{x}_i - \boldsymbol{\mu}\|^2}{\|\mathbf{x}_i - \boldsymbol{\mu}_m\|}$$

majorizes $f(\boldsymbol{\mu})$ at the current iterate $\boldsymbol{\mu}_m$ up to an irrelevant constant. Deduce Weiszfeld's algorithm [271]

$$\boldsymbol{\mu}_{m+1} = \frac{1}{\sum_{i=1}^p w_{mi}} \sum_{i=1}^p w_{mi} \mathbf{x}_i$$

with weights $w_{mi} = 1/\|\mathbf{x}_i - \boldsymbol{\mu}_m\|$. Comment on the relevance of this algorithm to a robust analogue of k -means clustering.

21. Show that the function $f(\boldsymbol{\mu})$ of Problem 20 is strictly convex whenever the points $\mathbf{x}_1, \dots, \mathbf{x}_p$ are not collinear.
22. Explain how Problem 31 of Chap. 4 delivers a quadratic majorization. Describe circumstances under which this majorization is apt to be poor.
23. In ℓ_1 regression show that the estimate of the parameter vector satisfies the equality

$$\sum_{i=1}^m \operatorname{sgn}[y_i - \mu_i(\boldsymbol{\theta})] \nabla \mu_i(\boldsymbol{\theta}) = \mathbf{0},$$

provided no residual $y_i - \mu_i(\boldsymbol{\theta}) = 0$ and the regression functions $\mu_i(\boldsymbol{\theta})$ are differentiable. What is the corresponding equality for the modified ℓ_1 criterion (8.20)?

24. Problem 12 of Chap. 7 deals with minimizing the quadratic function $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ subject to the constraints $x_i \geq 0$. If one drops the assumption that $\mathbf{A} = (a_{ij})$ is positive definite, it is still possible to devise an MM algorithm. Define matrices \mathbf{A}^+ and \mathbf{A}^- with entries $\max\{a_{ij}, 0\}$ and $-\min\{a_{ij}, 0\}$, respectively. Based on the fifth and sixth majorizations of Problem 2, derive the MM updates

$$x_{m+1,i} = x_{m,i} \left[\frac{-b_i + \sqrt{b_i^2 + 4(\mathbf{A}^+ \mathbf{x}_m)_i (\mathbf{A}^- \mathbf{x}_m)_i}}{2(\mathbf{A}^+ \mathbf{x}_m)_i} \right]$$

of Sha et al. [235]. All entries of the initial point \mathbf{x}_0 should be positive.

25. Show that the Bradley-Terry loglikelihood $f(\mathbf{r}) = \ln L(\mathbf{r})$ of Sect. 8.6 is concave under the reparameterization $r_i = e^{\theta_i}$.
26. In the Bradley-Terry model of Sect. 8.6, suppose we want to include the possibility of ties. One way of doing this is to write the probabilities of the three outcomes of i versus j as

$$\begin{aligned} \Pr(i \text{ wins}) &= \frac{r_i}{r_i + r_j + \theta \sqrt{r_i r_j}} \\ \Pr(i \text{ ties}) &= \frac{\theta \sqrt{r_i r_j}}{r_i + r_j + \theta \sqrt{r_i r_j}} \\ \Pr(i \text{ loses}) &= \frac{r_j}{r_i + r_j + \theta \sqrt{r_i r_j}}, \end{aligned}$$

where $\theta > 0$ is an additional parameter to be estimated. Let y_{ij} represent the number of times i beats j and t_{ij} the number of times i ties j . Prove that the loglikelihood of the data is

$$\begin{aligned} &L(\theta, \mathbf{r}) \\ &= \frac{1}{2} \sum_{i,j} \left(2y_{ij} \ln \frac{r_i}{r_i + r_j + \theta \sqrt{r_i r_j}} + t_{ij} \ln \frac{\theta \sqrt{r_i r_j}}{r_i + r_j + \theta \sqrt{r_i r_j}} \right). \end{aligned}$$

One way of maximizing $L(\theta, \mathbf{r})$ is to alternate between updating θ and \mathbf{r} . Both of these updates can be derived from the perspective of the MM algorithm. Two minorizations are now involved. The first proceeds using the convexity of $-\ln t$ just as in the text. This produces a function involving $-\sqrt{r_i r_j}$ terms. Demonstrate the minorization

$$-\sqrt{r_i r_j} \geq -\frac{r_i}{2} \sqrt{\frac{r_{mj}}{r_{mi}}} - \frac{r_j}{2} \sqrt{\frac{r_{mi}}{r_{mj}}},$$

and use it to minorize $L(\theta, \mathbf{r})$. Finally, determine \mathbf{r}_{m+1} for θ fixed at θ_m and θ_{m+1} for \mathbf{r} fixed at \mathbf{r}_{m+1} . The details are messy, but the overall strategy is straightforward.

27. In the linear logistic model of Sect. 8.7, it is possible to separate parameters and avoid matrix inversion altogether. In constructing a minorizing function, first prove the inequality

$$\begin{aligned} \ln[1 - \pi(\boldsymbol{\theta})] &= -\ln\left(1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}}\right) \\ &\geq -\ln\left(1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}\right) - \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}} - e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}}, \end{aligned}$$

with equality when $\boldsymbol{\theta} = \boldsymbol{\theta}_m$. This eliminates the log terms. Now apply the arithmetic-geometric mean inequality to the exponential functions $e^{\mathbf{x}_i^* \boldsymbol{\theta}}$ to separate parameters. Assuming that $\boldsymbol{\theta}$ has n components and that there are k observations, show that these maneuvers lead to the minorizing function

$$g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = -\frac{1}{n} \sum_{i=1}^k \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}} \sum_{j=1}^n e^{n x_{ij} (\theta_j - \theta_{mj})} + \sum_{i=1}^k y_i \mathbf{x}_i^* \boldsymbol{\theta}$$

up to a constant that does not depend on $\boldsymbol{\theta}$. Finally, prove that maximizing $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ consists in solving the transcendental equation

$$-\sum_{i=1}^k \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}_m} x_{ij} e^{-n x_{ij} \theta_{mj}}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}_m}} e^{n x_{ij} \theta_j} + \sum_{i=1}^k y_i x_{ij} = 0$$

for each j . This can be accomplished numerically.

28. Consider the general posynomial of n variables

$$f(\mathbf{x}) = \sum_{\boldsymbol{\alpha} \in S} c_{\boldsymbol{\alpha}} \prod_{i=1}^n x_i^{\alpha_i}$$

subject to the constraints $x_i > 0$ for each i . Here the index set $S \subset \mathbb{R}^n$ is finite, and the coefficients $c_{\boldsymbol{\alpha}}$ are positive. We can assume that at least one $\alpha_i > 0$ and at least one $\alpha_i < 0$ for every i . Otherwise, $f(\mathbf{x})$ can be reduced by sending x_i to ∞ or 0. Demonstrate that $f(\mathbf{x})$ is majorized by the sum

$$\begin{aligned} g(\mathbf{x} \mid \mathbf{x}_m) &= \sum_{i=1}^n g_i(x_i \mid \mathbf{x}_m) \\ g_i(x_i \mid \mathbf{x}_m) &= \sum_{\boldsymbol{\alpha} \in S} c_{\boldsymbol{\alpha}} \left(\prod_{j=1}^n x_{mj}^{\alpha_j} \right) \frac{|\alpha_i|}{\|\boldsymbol{\alpha}\|_1} \left(\frac{x_i}{x_{mi}} \right)^{\|\boldsymbol{\alpha}\|_1 \operatorname{sgn}(\alpha_i)}, \end{aligned}$$

where $\|\boldsymbol{\alpha}\|_1 = \sum_{j=1}^n |\alpha_j|$ and $\operatorname{sgn}(\alpha_i)$ is the sign function. To prove that the MM algorithm is well defined and produces iterates with positive entries, demonstrate that

$$\lim_{x_i \rightarrow \infty} g_i(x_i \mid \mathbf{x}_m) = \lim_{x_i \rightarrow 0} g_i(x_i \mid \mathbf{x}_m) = \infty.$$

Finally change variables by setting

$$\begin{aligned} y_i &= \ln x_i \\ h_i(y_i | \mathbf{x}_m) &= g_i(x_i | \mathbf{x}_m) \end{aligned}$$

for each i . Show that $h_i(y_i | \mathbf{x}_m)$ is strictly convex in y_i and therefore possesses a unique minimum point. The latter property carries over to the surrogate function $g_i(x_i | \mathbf{x}_m)$.

29. Devise MM algorithms based on Problem 28 to minimize the posynomials

$$\begin{aligned} f_1(\mathbf{x}) &= \frac{1}{x_1 x_2^2} + x_1 x_2^2 \\ f_2(\mathbf{x}) &= \frac{1}{x_1 x_2^2} + x_1 x_2. \end{aligned}$$

In the first case, demonstrate that the MM algorithm iterates according to

$$x_{m+1,1} = \sqrt[3]{\frac{x_{m1}^2}{x_{m2}^2}}, \quad x_{m+1,2} = \sqrt[3]{\frac{x_{m2}}{x_{m1}}}.$$

Furthermore, show that (a) $f_1(\mathbf{x})$ attains its minimum value of 2 whenever $x_1 x_2^2 = 1$, (b) the MM algorithm converges after a single iteration to the value 2, and (c) the converged point \mathbf{x}_1 depends on the initial point \mathbf{x}_0 . In the second case, demonstrate that the MM algorithm iterates according to

$$x_{m+1,1} = \sqrt[5]{\frac{x_{m1}^3}{x_{m2}^3}}, \quad x_{m+1,2} = \sqrt[5]{2 \frac{x_{m2}^2}{x_{m1}}}.$$

Furthermore, show that (a) the infimum of $f_2(\mathbf{x})$ is 0, (b) the MM algorithm satisfies the identities

$$x_{m1} x_{m2}^{3/2} = 2^{3/10}, \quad x_{m+1,2} = 2^{2/25} x_{m2}$$

for all $m \geq 2$, and (c) the minimum value 0 is attained asymptotically with x_{m1} tending to 0 and x_{m2} tending to ∞ .

30. A general posynomial of n variables can be represented as

$$h(\mathbf{y}) = \sum_{\alpha \in S} c_\alpha e^{\alpha^* \mathbf{y}}$$

in the parameterization $y_i = \ln x_i$. Here the index set $S \subset \mathbb{R}^n$ is finite and the coefficients c_α are positive. Show that $h(\mathbf{y})$ is strictly convex if and only if the power vectors $\alpha \in S$ span \mathbb{R}^n .

31. Demonstrate the minorization

$$\prod_{i=1}^n x_i^{\alpha_i} \geq \prod_{j=1}^n x_{mj}^{\alpha_j} \left(1 + \sum_{i=1}^n \alpha_i \ln x_i - \sum_{i=1}^n \alpha_i \ln x_{mi} \right).$$

How is this helpful in signomial programming [170]? As a concrete example, devise an MM algorithm to minimize the function

$$f(\mathbf{x}) = x_1^2 x_2^2 - 2x_1 x_2 x_3 x_4 + x_3^2 x_4^2 = (x_1 x_2 - x_3 x_4)^2.$$

Describe the solution as you vary the initial point.

32. Even more functions can be brought under the umbrella of signomial programming. For instance, majorization of the functions $-\ln f(\mathbf{x})$ and $\ln f(\mathbf{x})$ is possible for any posynomial

$$f(\mathbf{x}) = \sum_{\alpha \in S} c_{\alpha} \prod_{i=1}^n x_i^{\alpha_i}.$$

In the first case show that

$$-\ln f(\mathbf{x}) \leq -\sum_{\alpha \in S} \frac{d_{m\alpha}}{e_m} \left[\sum_{i=1}^n \alpha_i \ln x_i + \ln \left(\frac{c_{\alpha} e_m}{d_{m\alpha}} \right) \right] \quad (8.21)$$

holds for $d_{m\alpha} = c_{\alpha} \prod_{i=1}^n x_{mi}^{\alpha_i}$ and $e_m = \sum_{\alpha} d_{m\alpha}$. In the second case, show that

$$\ln f(\mathbf{x}) \leq \ln f(\mathbf{x}_m) + \frac{1}{f(\mathbf{x}_m)} [f(\mathbf{x}) - f(\mathbf{x}_m)].$$

This second majorization yields a posynomial, which can be majorized by the methods already described. Note that the coefficients c_{α} can be negative as well as positive in the second case [170].

33. Show that the loglikelihood (8.14) for the transmission tomography model is concave. State a necessary condition for strict concavity in terms of the number of pixels and the number of projections.

34. In the maximization phase of the MM algorithm for transmission tomography without a smoothing prior, demonstrate that the exact solution of the one-dimensional equation

$$\frac{\partial}{\partial \theta_j} g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = 0$$

exists and is positive when $\sum_i l_{ij} d_i > \sum_i l_{ij} y_i$. Why would this condition typically hold in practice?

35. Prove that the functions $\psi(r) = \sqrt{r^2 + \epsilon}$ and $\psi(r) = \ln[\cosh(r)]$ are even, strictly convex, infinitely differentiable, and asymptotic to $|r|$ as $|r| \rightarrow \infty$.
36. In positron emission tomography (PET), one seeks to estimate an object's Poisson emission intensities $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_p)^*$. Here p pixels are arranged in a 2-dimensional grid surrounded by an array of photon detectors [167, 265]. The observed data are coincidence counts (y_1, \dots, y_d) along d lines of flight connecting pairs of photon detectors. The loglikelihood under the PET model is

$$L(\boldsymbol{\lambda}) = \sum_i \left[y_i \ln \left(\sum_j e_{ij} \lambda_j \right) - \sum_j e_{ij} \lambda_j \right],$$

where the constants e_{ij} are determined from the geometry of the grid and the detectors. One can assume without loss of generality that $\sum_i e_{ij} = 1$. Use Jensen's inequality to derive the minorization

$$L(\boldsymbol{\lambda}) \geq \sum_i y_i \sum_j w_{nij} \ln \left(\frac{e_{ij} \lambda_j}{w_{nij}} \right) - \sum_i \sum_j e_{ij} \lambda_j = Q(\boldsymbol{\lambda} | \boldsymbol{\lambda}_n),$$

where $w_{nij} = e_{ij} \lambda_{nj} / (\sum_k e_{ik} \lambda_{nk})$. Show that the stationarity conditions for the surrogate function $Q(\boldsymbol{\lambda} | \boldsymbol{\lambda}_n)$ entail the MM updates

$$\lambda_{n+1,j} = \frac{\sum_i y_i w_{nij}}{\sum_i e_{ij}}.$$

To smooth the image, one can maximize the penalized loglikelihood

$$f(\boldsymbol{\lambda}) = L(\boldsymbol{\lambda}) - \frac{\mu}{2} \sum_{\{j,k\} \in \mathcal{N}} (\lambda_j - \lambda_k)^2,$$

where μ is a tuning constant, and \mathcal{N} is a neighborhood system that pairs spatially adjacent pixels. If we adopt the majorization

$$(\lambda_j - \lambda_k)^2 \leq \frac{1}{2}(2\lambda_j - \lambda_{nj} - \lambda_{nk})^2 + \frac{1}{2}(2\lambda_k - \lambda_{nj} - \lambda_{nk})^2,$$

then show that the revised stationarity conditions are

$$0 = \sum_i \left[\frac{y_i w_{nij}}{\lambda_j} - e_{ij} \right] - \mu \sum_{k: \{j,k\} \in \mathcal{N}_j} (2\lambda_j - \lambda_{nj} - \lambda_{nk}).$$

From these equations derive the smoothed MM updates

$$\lambda_{n+1,j} = \frac{-b_{nj} - \sqrt{b_{nj}^2 - 4a_j c_{nj}}}{2a_j},$$

where

$$\begin{aligned} a_j &= -2\mu \sum_{k \in \mathcal{N}_j} 1 \\ b_{nj} &= \mu \sum_{k \in \mathcal{N}'_j} (\lambda_{nj} + \lambda_{nk}) - 1 \\ c_{nj} &= \sum_i y_i w_{nij}. \end{aligned}$$

Note that $a_j < 0$, so we take the negative sign before the square root.

37. Program and test on real data either of the random multigraph algorithms (8.16) or (8.17). The internet is a rich source of random graph data.
38. The Dirichlet-multinomial distribution is used to model multivariate count data $\mathbf{x} = (x_1, \dots, x_d)^*$ that is too over-dispersed to be handled reliably by the multinomial distribution. Recall that the Dirichlet-multinomial distribution represents a mixture of multinomial distribution with a Dirichlet prior on the cell probabilities p_1, \dots, p_d . If $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^*$ denotes the Dirichlet parameter vector, Δ_d denotes the unit simplex in \mathbb{R}^d , and $|\boldsymbol{\alpha}| = \sum_{i=1}^d \alpha_i$ and $|\mathbf{x}| = \sum_{i=1}^d x_i$, then the discrete density of the Dirichlet-multinomial is

$$\begin{aligned} f(\mathbf{x} \mid \boldsymbol{\alpha}) &= \int_{\Delta_d} \binom{|\mathbf{x}|}{\mathbf{x}} \prod_{j=1}^d p_j^{x_j} \frac{\Gamma(|\boldsymbol{\alpha}|)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_d)} \prod_{j=1}^d p_j^{\alpha_j - 1} dp_1 \cdots dp_d \\ &= \binom{|\mathbf{x}|}{\mathbf{x}} \frac{\Gamma(\alpha_1 + x_1) \cdots \Gamma(\alpha_d + x_d)}{\Gamma(|\boldsymbol{\alpha}| + |\mathbf{x}|)} \frac{\Gamma(|\boldsymbol{\alpha}|)}{\Gamma(\alpha_1) \cdots \Gamma(\alpha_d)} \\ &= \binom{|\mathbf{x}|}{\mathbf{x}} \frac{\prod_{j=1}^d \alpha_j (\alpha_j + 1) \cdots (\alpha_j + x_j - 1)}{|\boldsymbol{\alpha}| (|\boldsymbol{\alpha}| + 1) \cdots (|\boldsymbol{\alpha}| + |\mathbf{x}| - 1)}. \end{aligned}$$

Devise an MM algorithm to maximize $\ln f(\mathbf{x} \mid \boldsymbol{\alpha})$ by applying the supporting hyperplane inequality to each term $-\ln(|\boldsymbol{\alpha}| + k)$ and Jensen's inequality to each term $\ln(\alpha_j + k)$. These minorizations are designed to separate parameters. Show that the MM updates are

$$\alpha_{m+1,j} = \alpha_{mj} \frac{\sum_k \frac{s_{jk}}{\alpha_{mj} + k}}{\sum_k \frac{r_k}{|\alpha_m| + k}}$$

for appropriate constants r_k and s_{jk} . This problem and similar problems are treated in the reference [282].

39. In the dictionary model of motif finding [227], a DNA sequence is viewed as a concatenation of words independently drawn from a dictionary having the four letters A, C, G, and T. The words of the

dictionary of length k have collective probability q_k . The EM algorithm offers one method of estimating the q_k . Omitting many details, the EM algorithm maximizes the function

$$Q(\mathbf{q} \mid \mathbf{q}_m) = \sum_{k=1}^l c_{mk} \ln q_k - \ln \left(\sum_{k=1}^l k q_k \right).$$

Here the constants c_{mk} are positive, l is the maximum word length, and maximization is performed subject to the constraints $q_k \geq 0$ for $k = 1, \dots, l$ and $\sum_{k=1}^l q_k = 1$. Because this problem can not be solved in closed form, it is convenient to follow the EM minorization with a second minorization based on the inequality

$$\ln x \leq \ln y + x/y - 1. \quad (8.22)$$

Application of inequality (8.22) produces the minorizing function

$$h(\mathbf{q} \mid \mathbf{q}_m) = \sum_{k=1}^l c_{mk} \ln q_k - \ln \left(\sum_{k=1}^l k q_{mk} \right) - d_m \sum_{k=1}^l k q_k + 1$$

with $d_m = 1/(\sum_{k=1}^l k q_{mk})$.

- (a) Show that the function $h(\mathbf{q} \mid \mathbf{q}_m)$ minorizes $Q(\mathbf{q} \mid \mathbf{q}_m)$.
- (b) Maximize $h(\mathbf{q} \mid \mathbf{q}_m)$ using the method of Lagrange multipliers. At the current iteration, show that the solution has components

$$q_k = \frac{c_{mk}}{d_m k - \lambda}$$

for an unknown Lagrange multiplier λ .

- (c) Using the constraints, prove that λ exists and is unique.
- (d) Describe a reliable method for computing λ .
- (e) As an alternative to the exact method, construct a quadratic approximation to $h(\mathbf{q} \mid \mathbf{q}_m)$ near \mathbf{q}_m of the form

$$\frac{1}{2}(\mathbf{q} - \mathbf{q}_m)^* \mathbf{A}(\mathbf{q} - \mathbf{q}_m) + \mathbf{b}^*(\mathbf{q} - \mathbf{q}_m) + c.$$

In particular, what are \mathbf{A} and \mathbf{b} ?

- (f) Show that the quadratic approximation has maximum

$$\mathbf{q} = \mathbf{q}_m - \mathbf{A}^{-1} \left(\mathbf{b} - \frac{\mathbf{1}^* \mathbf{A}^{-1} \mathbf{b}}{\mathbf{1}^* \mathbf{A}^{-1} \mathbf{1}} \mathbf{1} \right) \quad (8.23)$$

subject to the constraint $\sum_{k=1}^l (q_k - q_{mk}) = 0$.

- (g) In the dictionary model, demonstrate that the solution (8.23) takes the form

$$q_j = q_{mj} + \frac{(q_{mj})^2}{c_{mj}} \left[\frac{c_{mj}}{q_{mj}} - d_{mj} - \frac{1 - \sum_{k=1}^l d_{mk} \frac{(q_{mk})^2}{c_{mk}}}{\sum_{k=1}^l \frac{(q_{mk})^2}{c_{mk}}} \right]$$

for the j th component of q .

- (h) Point out two potential pitfalls of this particular solution in conjunction with maximizing $h(\mathbf{q} \mid \mathbf{q}_m)$.

40. In the balanced ANOVA model with two factors, we estimate the parameter vector $\boldsymbol{\theta} = (\mu, \boldsymbol{\alpha}^*, \boldsymbol{\beta}^*)^*$ by minimizing the sum of squares

$$f(\boldsymbol{\theta}) = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K w_{ijk} (y_{ijk} - \mu - \alpha_i - \beta_j)^2$$

with all weights $w_{ijk} = 1$. If some of the observations y_{ijk} are missing, then we take the corresponding weights to be 0. The missing observations are now irrelevant, but it is possible to replace each one by its predicted value

$$\hat{y}_{ijk} = \mu + \alpha_i + \beta_j$$

given the current parameter values. If there are missing observations, de Leeuw [59] notes that

$$g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (z_{ijk} - \mu - \alpha_i - \beta_j)^2$$

majorizes $f(\boldsymbol{\theta})$ provided we define

$$z_{ijk} = \begin{cases} y_{ijk} & \text{for a regular observation} \\ \hat{y}_{ijk} & \text{for a missing observation.} \end{cases}$$

Prove this fact and calculate the MM update of $\boldsymbol{\theta}$, assuming the sum to 0 constraints $\sum_{i=1}^I \alpha_i = \sum_{j=1}^J \beta_j = 0$.

41. Consider the weighted sum of squares criterion

$$f(\boldsymbol{\theta}) = \sum_i w_i [y_i - \mu_i(\boldsymbol{\theta})]^2$$

with weights w_i drawn from the unit interval $[0, 1]$. Show that the function

$$\begin{aligned} g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) &= \sum_i [w_i y_i + (1 - w_i) \mu_i(\boldsymbol{\theta}_m) - \mu_i(\boldsymbol{\theta})]^2 + c_m \\ c_m &= \sum_i w_i (1 - w_i) [y_i - \mu_i(\boldsymbol{\theta}_m)]^2 \end{aligned}$$

majorizes $f(\boldsymbol{\theta})$ [121, 153, 242]. This majorization converts a weighted least squares problem into an ordinary least squares problem. (Hints: Write $[y_i - \mu_i(\boldsymbol{\theta})]^2 = [y_i - \mu_i(\boldsymbol{\theta}_m) + \mu_i(\boldsymbol{\theta}_m) - \mu_i(\boldsymbol{\theta})]^2$, expand, majorize $w_i[\mu_i(\boldsymbol{\theta}_m) - \mu_i(\boldsymbol{\theta})]^2$, and complete the square.)

42. Luce's model [181, 182] is a convenient scheme for ranking items such as candidates in an election, consumer goods in a certain category, or academic departments in a reputational survey. Some people will be too lazy or uncertain to rank each and every item, preferring to focus on just their top choices. How can we use this form of limited voting to rank the entire set of items? A partial ranking by a person is a sequence of random choices X_1, \dots, X_l , with X_1 the highest ranked item, X_2 the second highest ranked item, and so forth. If there are r items, then the number l of items ranked may be strictly less than r ; $l = 1$ is a distinct possibility. The data arrive at our doorstep as a random sample of s independent partial rankings, which we must integrate in some coherent fashion. One possibility is to adopt multinomial sampling without replacement. This differs from ordinary multinomial sampling in that once an item is chosen, it cannot be chosen again. However, remaining items are selected with the conditional probabilities dictated by the original sampling probabilities. Show that the likelihood under Luce's model reduces to

$$\begin{aligned} & \prod_{i=1}^s \Pr(X_{i1} = x_{i1}, \dots, X_{il_i} = x_{il_i}) \\ &= \prod_{i=1}^s p_{x_{i1}} \prod_{j=1}^{l_i-1} \frac{p_{x_{i,j+1}}}{\sum_{k \notin \{x_{i1}, \dots, x_{ij}\}} p_k}, \end{aligned}$$

where x_{ij} is the j th choice out of l_i choices for person i and p_k is the multinomial probability assigned to item k . If we can estimate the p_k , then we can rank the items accordingly. The item with largest estimated probability is ranked first and so on.

The model has the added virtue of leading to straightforward estimation by the MM algorithm. Use the supporting hyperplane inequality

$$-\ln t \geq -\ln t_m - \frac{1}{t_m}(t - t_m).$$

to generate the minorization

$$Q(\mathbf{p} \mid \mathbf{p}_m) = \sum_{i=1}^s \sum_{j=1}^{l_i} \ln p_{x_{ij}} - \sum_{i=1}^s \sum_{j=1}^{l_i-1} w_{ij} \sum_{k \notin \{x_{i1}, \dots, x_{ij}\}} p_k$$

of the loglikelihood up to an irrelevant constant. Specify the positive weights w_{ij} and derive the maximization step of the MM algorithm.

Show that your update has the intuitive interpretation of equating the expected number of choices of item k to the observed number of choices of item k across all voters. Finally, generalize the model so that person i 's choices are limited to a subset S_i of the items. For instance, in rating academic departments, some people may only feel competent to rank those departments in their state or region. What form does the MM algorithm take in this setting?

9

The EM Algorithm

9.1 Introduction

Maximum likelihood is the dominant form of estimation in applied statistics. Because closed-form solutions to likelihood equations are the exception rather than the rule, numerical methods for finding maximum likelihood estimates are of paramount importance. In this chapter we study maximum likelihood estimation by the EM algorithm [65, 179, 191], a special case of the MM algorithm. At the heart of every EM algorithm is some notion of missing data. Data can be missing in the ordinary sense of a failure to record certain observations on certain cases. Data can also be missing in a theoretical sense. We can think of the E (expectation) step of the algorithm as filling in the missing data. This action replaces the log-likelihood of the observed data by a minorizing function. This surrogate function is then maximized in the M step. Because the surrogate function is usually much simpler than the likelihood, we can often solve the M step analytically. The price we pay for this simplification is that the EM algorithm is iterative. Reconstructing the missing data is bound to be slightly wrong if the parameters do not already equal their maximum likelihood estimates.

One of the advantages of the EM algorithm is its numerical stability. As an MM algorithm, any EM algorithm leads to a steady increase in the likelihood of the observed data. Thus, the EM algorithm avoids wildly overshooting or undershooting the maximum of the likelihood along its current direction of search. Besides this desirable feature, the EM handles

parameter constraints gracefully. Constraint satisfaction is by definition built into the solution of the M step. In contrast, competing methods of maximization must incorporate special techniques to cope with parameter constraints. The EM shares some of the negative features of the more general MM algorithm. For example, the EM algorithm often converges at an excruciatingly slow rate in a neighborhood of the maximum point. This rate directly reflects the amount of missing data in a problem. In the absence of concavity, there is also no guarantee that the EM algorithm will converge to the global maximum. The global maximum can usually be reached by starting the parameters at good but suboptimal estimates such as method-of-moments estimates or by choosing multiple random starting points.

9.2 Definition of the EM Algorithm

A sharp distinction is drawn in the EM algorithm between the observed, incomplete data \mathbf{y} and the unobserved, complete data \mathbf{x} of a statistical experiment [65, 179, 251]. Some function $t(\mathbf{x}) = \mathbf{y}$ collapses \mathbf{x} onto \mathbf{y} . For instance, if we represent \mathbf{x} as (\mathbf{y}, \mathbf{z}) , with \mathbf{z} as the missing data, then t is simply projection onto the \mathbf{y} -component of \mathbf{x} . It should be stressed that the missing data can consist of more than just observations missing in the ordinary sense. In fact, the definition of \mathbf{x} is left up to the intuition and cleverness of the statistician. The general idea is to choose \mathbf{x} so that maximum likelihood estimation becomes trivial for the complete data.

The complete data are assumed to have a probability density $f(\mathbf{x} | \boldsymbol{\theta})$ that is a function of a parameter vector $\boldsymbol{\theta}$ as well as of \mathbf{x} . In the E step of the EM algorithm, we calculate the conditional expectation

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n) = E[\ln f(\mathbf{X} | \boldsymbol{\theta}) | \mathbf{Y} = \mathbf{y}, \boldsymbol{\theta}_n].$$

Here $\boldsymbol{\theta}_n$ is the current estimated value of $\boldsymbol{\theta}$, upper case letters indicate random vectors, and lower case letters indicate corresponding realizations of these random vectors. In the M step, we maximize $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$ with respect to $\boldsymbol{\theta}$. This yields the new parameter estimate $\boldsymbol{\theta}_{n+1}$, and we repeat this two-step process until convergence occurs. Note that $\boldsymbol{\theta}$ and $\boldsymbol{\theta}_n$ play fundamentally different roles in $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$.

If $\ln g(\mathbf{y} | \boldsymbol{\theta})$ denotes the loglikelihood of the observed data, then the EM algorithm enjoys the ascent property

$$\ln g(\mathbf{y} | \boldsymbol{\theta}_{n+1}) \geq \ln g(\mathbf{y} | \boldsymbol{\theta}_n).$$

Proof of this assertion unfortunately involves measure theory, so some readers may want to take it on faith and skip the rest of this section. A necessary preliminary is the following well-known inequality from statistics.

Proposition 9.2.1 (Information Inequality) *Let h and k be probability densities with respect to a measure μ . Suppose $h > 0$ and $k > 0$ almost everywhere relative to μ . If E_h denotes expectation with respect to the probability measure $h d\mu$, then $E_h(\ln h) \geq E_h(\ln k)$, with equality if and only if $h = k$ almost everywhere relative to μ .*

Proof: Because $-\ln(w)$ is a strictly convex function on $(0, \infty)$, Proposition 6.6.1 applied to the random variable k/h yields

$$\begin{aligned} E_h(\ln h) - E_h(\ln k) &= E_h\left(-\ln \frac{k}{h}\right) \\ &\geq -\ln E_h\left(\frac{k}{h}\right) \\ &= -\ln \int \frac{k}{h} h d\mu \\ &= -\ln \int k d\mu \\ &= 0. \end{aligned}$$

Equality holds if and only if $k/h = E_h(k/h)$ almost everywhere relative to μ . This necessary and sufficient condition is equivalent to $h = k$ since $E_h(k/h) = 1$. ■

To prove the ascent property of the EM algorithm, it suffices to demonstrate the minorization inequality

$$\ln g(\mathbf{y} | \boldsymbol{\theta}) \geq Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n) + \ln g(\mathbf{y} | \boldsymbol{\theta}_n) - Q(\boldsymbol{\theta}_n | \boldsymbol{\theta}_n),$$

where $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n) = E[\ln f(\mathbf{X} | \boldsymbol{\theta}) | \mathbf{Y} = \mathbf{y}, \boldsymbol{\theta}_n]$. With this end in mind, note that both $f(\mathbf{x} | \boldsymbol{\theta})/g(\mathbf{y} | \boldsymbol{\theta})$ and $f(\mathbf{x} | \boldsymbol{\theta}_n)/g(\mathbf{y} | \boldsymbol{\theta}_n)$ are conditional densities of \mathbf{X} on $\{\mathbf{x} : t(\mathbf{x}) = \mathbf{y}\}$ with respect to some measure $\mu_{\mathbf{y}}$. The information inequality now indicates that

$$\begin{aligned} Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n) - \ln g(\mathbf{y} | \boldsymbol{\theta}) &= E\left(\ln \left[\frac{f(\mathbf{X} | \boldsymbol{\theta})}{g(\mathbf{Y} | \boldsymbol{\theta})}\right] \middle| \mathbf{Y} = \mathbf{y}, \boldsymbol{\theta}_n\right) \\ &\leq E\left(\ln \left[\frac{f(\mathbf{X} | \boldsymbol{\theta}_n)}{g(\mathbf{Y} | \boldsymbol{\theta}_n)}\right] \middle| \mathbf{Y} = \mathbf{y}, \boldsymbol{\theta}_n\right) \\ &= Q(\boldsymbol{\theta}_n | \boldsymbol{\theta}_n) - \ln g(\mathbf{y} | \boldsymbol{\theta}_n). \end{aligned}$$

Maximizing $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$ therefore drives $\ln g(\mathbf{y} | \boldsymbol{\theta})$ uphill. The ascent inequality is strict whenever the conditional density $f(\mathbf{x} | \boldsymbol{\theta})/g(\mathbf{y} | \boldsymbol{\theta})$ differs at the parameter points $\boldsymbol{\theta}_n$ and $\boldsymbol{\theta}_{n+1}$ or

$$Q(\boldsymbol{\theta}_{n+1} | \boldsymbol{\theta}_n) > Q(\boldsymbol{\theta}_n | \boldsymbol{\theta}_n).$$

The preceding proof is a little vague as to the meaning of the conditional density $f(\mathbf{x} | \boldsymbol{\theta})/g(\mathbf{y} | \boldsymbol{\theta})$ and its associated measure $\mu_{\mathbf{y}}$. Commonly the

complete data decomposes as $\mathbf{x} = (\mathbf{y}, \mathbf{z})$, where \mathbf{z} is considered the missing data and $t(\mathbf{y}, \mathbf{z}) = \mathbf{y}$ is projection onto the observed data. Suppose (\mathbf{y}, \mathbf{z}) has joint density $f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta})$ relative to a product measure $\omega \times \mu(\mathbf{y}, \mathbf{z})$; ω and μ are typically Lebesgue measure or counting measure. In this framework, we define $g(\mathbf{y} | \boldsymbol{\theta}) = \int f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) d\mu(\mathbf{z})$ and set $\mu_{\mathbf{y}} = \mu$. The function $g(\mathbf{y} | \boldsymbol{\theta})$ serves as a density relative to ω . To check that these definitions make sense, it suffices to prove that $\int h(\mathbf{y}, \mathbf{z}) f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) / g(\mathbf{y} | \boldsymbol{\theta}) d\mu(\mathbf{z})$ is a version of the conditional expectation $E[h(\mathbf{Y}, \mathbf{Z}) | \mathbf{Y} = \mathbf{y}]$ for every well-behaved function $h(\mathbf{y}, \mathbf{z})$. This assertion can be verified by showing

$$E\{1_S(\mathbf{Y}) E[h(\mathbf{Y}, \mathbf{Z}) | \mathbf{Y}]\} = E[1_S(\mathbf{Y}) h(\mathbf{Y}, \mathbf{Z})]$$

for every measurable set S . With

$$E[h(\mathbf{Y}, \mathbf{Z}) | \mathbf{Y} = \mathbf{y}] = \int h(\mathbf{y}, \mathbf{z}) \frac{f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta})}{g(\mathbf{y} | \boldsymbol{\theta})} d\mu(\mathbf{z}),$$

we calculate

$$\begin{aligned} & E\{1_S(\mathbf{Y}) E[h(\mathbf{Y}, \mathbf{Z}) | \mathbf{Y} = \mathbf{y}]\} \\ &= \int_S \int h(\mathbf{y}, \mathbf{z}) \frac{f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta})}{g(\mathbf{y} | \boldsymbol{\theta})} d\mu(\mathbf{z}) g(\mathbf{y} | \boldsymbol{\theta}) d\omega(\mathbf{y}) \\ &= \int_S \int h(\mathbf{y}, \mathbf{z}) f(\mathbf{y}, \mathbf{z} | \boldsymbol{\theta}) d\mu(\mathbf{z}) d\omega(\mathbf{y}) \\ &= E[1_S(\mathbf{Y}) h(\mathbf{Y}, \mathbf{Z})]. \end{aligned}$$

Hence in this situation, $f(\mathbf{x} | \boldsymbol{\theta}) / g(\mathbf{y} | \boldsymbol{\theta})$ is indeed the conditional density of \mathbf{X} given $\mathbf{Y} = \mathbf{y}$.

9.3 Missing Data in the Ordinary Sense

The most common application of the EM algorithm is to data missing in the ordinary sense. For example, Problem 40 of Chap. 8 considers a balanced ANOVA model with two factors. Missing observations in this setting break the symmetry that permits explicit solution of the likelihood equations. Thus, there is ample incentive for filling in the missing observations. If the observations follow an exponential model, and missing data are missing completely at random, then the EM algorithm replaces the sufficient statistic of each missing observation by its expected value.

The density of a random variable Y from an exponential family can be written as

$$f(y | \boldsymbol{\theta}) = g(y) e^{\beta(\boldsymbol{\theta}) + h(y)^* \boldsymbol{\gamma}(\boldsymbol{\theta})} \quad (9.1)$$

relative to some measure ν [73, 218]. The normal, Poisson, binomial, negative binomial, gamma, beta, and multinomial families are prime examples

of exponential families. The function $h(y)$ in equation (9.1) is the sufficient statistic. The maximum likelihood estimate of the parameter vector $\boldsymbol{\theta}$ depends on an observation y only through $h(y)$. Predictors of y are incorporated into the functions $\beta(\boldsymbol{\theta})$ and $\gamma(\boldsymbol{\theta})$.

To fill in a missing observation y , we take the ordinary expectation

$$E[\ln f(Y | \boldsymbol{\theta}) | \boldsymbol{\theta}_n] = E[\ln g(Y) | \boldsymbol{\theta}_n] + \beta(\boldsymbol{\theta}) + E[h(Y) | \boldsymbol{\theta}_n]^* \gamma(\boldsymbol{\theta})$$

of the complete data loglikelihood. This function is added to the loglikelihood of the regular observations y_1, \dots, y_m to generate the surrogate function $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$. For example, if a typical observation is normally distributed with mean $\mu(\boldsymbol{\alpha})$ and variance σ^2 , then $\boldsymbol{\theta}$ is the vector $(\boldsymbol{\alpha}^*, \sigma^2)^*$ and

$$\begin{aligned} E[\ln f(Y | \boldsymbol{\theta}) | \boldsymbol{\theta}_n] &= \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} E\{[Y - \mu(\boldsymbol{\alpha})]^2 | \boldsymbol{\theta}_n\} \\ &= \ln \frac{1}{\sqrt{2\pi\sigma^2}} - \frac{1}{2\sigma^2} \{\sigma_n^2 + [\mu(\boldsymbol{\alpha}_n) - \mu(\boldsymbol{\alpha})]^2\}. \end{aligned}$$

Once we have filled in the missing data, we can estimate $\boldsymbol{\alpha}$ without reference to σ^2 . This is accomplished by adding each square $[\mu_i(\boldsymbol{\alpha}_n) - \mu_i(\boldsymbol{\alpha})]^2$ corresponding to a missing observation y_i to the sum of squares for the actual observations and then minimizing the entire sum over $\boldsymbol{\alpha}$. In classical models such as balanced ANOVA, the M step is exact. Once the iterative limit $\lim_{n \rightarrow \infty} \boldsymbol{\alpha}_n = \hat{\boldsymbol{\alpha}}$ is reached, we can estimate σ^2 in one step by the formula

$$\hat{\sigma}^2 = \frac{1}{m} \sum_{i=1}^m [y_i - \mu_i(\hat{\boldsymbol{\alpha}})]^2$$

using only the observed y_i . The reader is urged to work Problem 40 of Chap. 8 to see the whole process in action.

9.4 Allele Frequency Estimation

It is instructive to compare the EM and MM algorithms on identical problems. Even when the two algorithms specify the same iteration scheme, the differences in deriving the algorithms are illuminating. Consider the ABO allele frequency estimation problem of Sect. 8.4. From the EM perspective, the complete data \boldsymbol{x} are genotype counts rather than phenotype counts \boldsymbol{y} . In passing from the complete data to the observed data, nature collapses genotypes A/A and A/O into phenotype A and genotypes B/B and B/O into phenotype B . In view of the Hardy-Weinberg equilibrium law, the complete data multinomial loglikelihood becomes

$$\ln f(X | \boldsymbol{p}) = n_{A/A} \ln p_A^2 + n_{A/O} \ln(2p_A p_O) + n_{B/B} \ln p_B^2$$

$$\begin{aligned}
& + n_{B/O} \ln(2p_B p_O) + n_{AB} \ln(2p_A p_B) + n_O \ln p_O^2 \\
& + \ln \binom{n}{n_{A/A}, n_{A/O}, n_{B/B}, n_{B/O}, n_{AB}, n_O}. \quad (9.2)
\end{aligned}$$

In the E step of the EM algorithm we take the expectation of $\ln f(X | \mathbf{p})$ conditional on the observed counts n_A, n_B, n_{AB} , and n_O and the current parameter vector $\mathbf{p}_m = (p_{mA}, p_{mB}, p_{mO})^*$. This action yields the surrogate function

$$\begin{aligned}
Q(\mathbf{p} | \mathbf{p}_m) &= \mathbb{E}(n_{A/A} | Y, \mathbf{p}_m) \ln p_A^2 + \mathbb{E}(n_{A/O} | Y, \mathbf{p}_m) \ln(2p_A p_O) \\
&+ \mathbb{E}(n_{B/B} | Y, \mathbf{p}_m) \ln p_B^2 + \mathbb{E}(n_{B/O} | Y, \mathbf{p}_m) \ln(2p_B p_O) \\
&+ \mathbb{E}(n_{AB} | Y, \mathbf{p}_m) \ln(2p_A p_B) + \mathbb{E}(n_O | Y, \mathbf{p}_m) \ln p_O^2 \\
&+ \mathbb{E} \left[\ln \binom{n}{n_{A/A}, n_{A/O}, n_{B/B}, n_{B/O}, n_{AB}, n_O} \middle| Y, \mathbf{p}_m \right].
\end{aligned}$$

It is obvious that

$$\begin{aligned}
\mathbb{E}(n_{AB} | Y, \mathbf{p}_m) &= n_{AB} \\
\mathbb{E}(n_O | Y, \mathbf{p}_m) &= n_O.
\end{aligned}$$

Application of Bayes' rule gives

$$\begin{aligned}
n_{mA/A} &= \mathbb{E}(n_{A/A} | Y, \mathbf{p}_m) \\
&= n_A \frac{p_{mA}^2}{p_{mA}^2 + 2p_{mA} p_{mO}} \\
n_{mA/O} &= \mathbb{E}(n_{A/O} | Y, \mathbf{p}_m) \\
&= n_A \frac{2p_{mA} p_{mO}}{p_{mA}^2 + 2p_{mA} p_{mO}}.
\end{aligned}$$

The conditional expectations $n_{mB/B}$ and $n_{mB/O}$ reduce to similar expressions. Hence, the surrogate function $Q(\mathbf{p} | \mathbf{p}_m)$ derived from the complete data likelihood matches the surrogate function of the MM algorithm up to a constant, and the maximization step proceeds as described earlier. One of the advantages of the EM derivation is that it explicitly reveals the nature of the conditional expectations $n_{mA/A}$, $n_{mA/O}$, $n_{mB/B}$, and $n_{mB/O}$.

9.5 Clustering by EM

The k-means clustering algorithm discussed in Example 7.2.3 makes hard choices in cluster assignment. The alternative of soft choices is possible with admixture models [192, 259]. An admixture probability density $h(\mathbf{y})$ can be written as a convex combination

$$h(\mathbf{y}) = \sum_{j=1}^k \pi_j h_j(\mathbf{y}), \quad (9.3)$$

where the π_j are nonnegative probabilities that sum to 1 and $h_j(\mathbf{y})$ is the probability density of group j . According to Bayes' rule, the posterior probability that an observation \mathbf{y} belongs to group j equals the ratio

$$\frac{\pi_j h_j(\mathbf{y})}{\sum_{i=1}^k \pi_i h_i(\mathbf{y})}. \tag{9.4}$$

If hard assignment is necessary, then the rational procedure is to assign \mathbf{y} to the group with highest posterior probability.

Suppose the observations $\mathbf{y}_1, \dots, \mathbf{y}_m$ represent a random sample from the admixture density (9.3). In practice we want to estimate the admixture proportions and whatever further parameters θ characterize the densities $h_j(\mathbf{y} \mid \theta)$. The EM algorithm is natural in this context with group membership as the missing data. If we let z_{ij} be an indicator specifying whether observation \mathbf{y}_i comes from group j , then the complete data loglikelihood amounts to

$$\sum_{i=1}^m \sum_{j=1}^k z_{ij} [\ln \pi_j + \ln h_j(\mathbf{y}_i \mid \theta)].$$

To find the surrogate function, we must find the conditional expectation w_{ij} of z_{ij} . But this reduces to the Bayes' rule (9.4) with θ fixed at θ_n and π fixed at π_n , where as usual n indicates iteration number. Note that the property $\sum_{j=1}^k z_{ij} = 1$ entails the property $\sum_{j=1}^k w_{ij} = 1$.

Fortunately, the E step of the EM algorithm separates the π parameters from the θ parameters. The problem of maximizing

$$\sum_{j=1}^k c_j \ln \pi_j$$

with $c_j = \sum_{i=1}^m w_{ij}$ should be familiar by now. Since $\sum_{j=1}^k c_j = m$, Example (1.4.2) shows that $\pi_{n+1,j} = \frac{c_j}{m}$.

We now undertake estimation of the remaining parameters assuming the groups are normally distributed with a common variance matrix Ω but different mean vectors μ_1, \dots, μ_k . The pertinent part of the surrogate function is

$$\begin{aligned} & \sum_{i=1}^m \sum_{j=1}^k w_{ij} \left[-\frac{1}{2} \ln \det \Omega - \frac{1}{2} (\mathbf{y}_i - \mu_j)^* \Omega^{-1} (\mathbf{y}_i - \mu_j) \right] \\ &= -\frac{m}{2} \ln \det \Omega - \frac{1}{2} \sum_{j=1}^k \sum_{i=1}^m w_{ij} (\mathbf{y}_i - \mu_j)^* \Omega^{-1} (\mathbf{y}_i - \mu_j) \\ &= -\frac{m}{2} \ln \det \Omega - \frac{1}{2} \operatorname{tr} \left[\Omega^{-1} \sum_{j=1}^k \sum_{i=1}^m w_{ij} (\mathbf{y}_i - \mu_j) (\mathbf{y}_i - \mu_j)^* \right]. \tag{9.5} \end{aligned}$$

Differentiating the surrogate (9.5) with respect to $\boldsymbol{\mu}_j$ gives the equation

$$\sum_{i=1}^m w_{ij} \boldsymbol{\Omega}^{-1} (\mathbf{y}_i - \boldsymbol{\mu}_j) = \mathbf{0}$$

with solution

$$\boldsymbol{\mu}_{n+1,j} = \frac{1}{\sum_{i=1}^m w_{ij}} \sum_{i=1}^m w_{ij} \mathbf{y}_i.$$

Maximization of the surrogate (9.5) with respect to $\boldsymbol{\Omega}$ can be rephrased as maximization of

$$-\frac{m}{2} \ln \det \boldsymbol{\Omega} - \frac{1}{2} \text{tr}(\boldsymbol{\Omega}^{-1} \mathbf{M})$$

for the choice

$$\mathbf{M} = \sum_{j=1}^k \sum_{i=1}^m w_{ij} (\mathbf{y}_i - \boldsymbol{\mu}_{n+1,j})(\mathbf{y}_i - \boldsymbol{\mu}_{n+1,j})^*.$$

Abstractly this is just the problem we faced in Example 6.5.7. Inspection of the arguments there shows that

$$\boldsymbol{\Omega}_{n+1} = \frac{1}{m} \mathbf{M}. \quad (9.6)$$

There is no guarantee of a unique mode in this model. Fortunately, k-means clustering generates good starting values for the parameters. The cluster centers provide the group means. If we set w_{ij} equal to 1 or 0 depending on whether observation i belongs to cluster j or not, then the matrix (9.6) serves as an initial guess of the common variance matrix. The initial admixture proportion π_j can be taken to be the proportion of the observations assigned to cluster j .

9.6 Transmission Tomography

The EM and MM algorithms for transmission tomography differ. The MM algorithm is easier to derive and computationally more efficient. In other examples, the opposite is true.

In the transmission tomography example of Sect. 8.10, it is natural to view the missing data as the number of photons X_{ij} entering each pixel j along each projection line i . These random variables supplemented by the observations Y_i constitute the complete data. If projection line i does not intersect pixel j , then $X_{ij} = 0$. Although X_{ij} and $X_{i'j'}$ are not independent,

the collection $\{X_{ij}\}_j$ indexed by projection i is independent of the collection $\{X_{i'j}\}_j$ indexed by another projection i' . This allows us to work projection by projection in writing the complete data likelihood. We will therefore temporarily drop the projection subscript i and relabel pixels, starting with pixel 1 adjacent to the source and ending with pixel $m - 1$ adjacent to the detector. In this notation X_1 is the number of photons leaving the source, X_j is the number of photons entering pixel j , and $X_m = Y$ is the number of photons detected.

By assumption X_1 follows a Poisson distribution with mean d . Conditional on X_1, \dots, X_j , the random variable X_{j+1} is binomially distributed with X_j trials and success probability $e^{-l_j\theta_j}$. In other words, each of the X_j photons entering pixel j behaves independently and has a chance $e^{-l_j\theta_j}$ of avoiding attenuation in pixel j . It follows that the complete data loglikelihood for the current projection is

$$\begin{aligned}
 & -d + X_1 \ln d - \ln X_1! \\
 & + \sum_{j=1}^{m-1} \left[\ln \binom{X_j}{X_{j+1}} + X_{j+1} \ln e^{-l_j\theta_j} + (X_j - X_{j+1}) \ln(1 - e^{-l_j\theta_j}) \right].
 \end{aligned} \tag{9.7}$$

To perform the E step of the EM algorithm, we merely need to compute the conditional expectations $E(X_j \mid X_m = y, \boldsymbol{\theta})$ for $1 \leq j \leq m$. The conditional expectations of other terms such as $\ln \binom{X_j}{X_{j+1}}$ appearing in (9.7) are irrelevant in the subsequent M step.

Reasoning as above, we infer that the unconditional mean of X_j is

$$\mu_j = E(X_j) = de^{-\sum_{k=1}^{j-1} l_k\theta_k}$$

and that the distribution of X_m conditional on X_j is binomial with X_j trials and success probability

$$\frac{\mu_m}{\mu_j} = e^{-\sum_{k=j}^{m-1} l_k\theta_k}.$$

In view of our remarks about random thinning in Chap. 8, the joint probability density of X_j and X_m therefore reduces to

$$\Pr(X_j = x_j, X_m = x_m) = e^{-\mu_j} \frac{\mu_j^{x_j}}{x_j!} \binom{x_j}{x_m} \left(\frac{\mu_m}{\mu_j}\right)^{x_m} \left(1 - \frac{\mu_m}{\mu_j}\right)^{x_j - x_m},$$

and the conditional probability density of X_j given X_m becomes

$$\begin{aligned}
 \Pr(X_j = x_j \mid X_m = x_m) &= \frac{e^{-\mu_j} \frac{\mu_j^{x_j}}{x_j!} \binom{x_j}{x_m} \left(\frac{\mu_m}{\mu_j}\right)^{x_m} \left(1 - \frac{\mu_m}{\mu_j}\right)^{x_j - x_m}}{e^{-\mu_m} \frac{\mu_m^{x_m}}{x_m!}} \\
 &= e^{-(\mu_j - \mu_m)} \frac{(\mu_j - \mu_m)^{x_j - x_m}}{(x_j - x_m)!}.
 \end{aligned}$$

In other words, conditional on X_m , the difference $X_j - X_m$ follows a Poisson distribution with mean $\mu_j - \mu_m$. This implies in particular that

$$\begin{aligned} E(X_j | X_m) &= E(X_j - X_m | X_m) + X_m \\ &= \mu_j - \mu_m + X_m. \end{aligned}$$

Reverting to our previous notation, it is now possible to assemble the surrogate function $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$ of the E step. Define

$$\begin{aligned} M_{ij} &= d_i(e^{-\sum_{k \in S_{ij}} l_{ik}\theta_{nk}} - e^{-\sum_k l_{ik}\theta_{nk}}) + y_i \\ N_{ij} &= d_i(e^{-\sum_{k \in S_{ij} \cup \{j\}} l_{ik}\theta_{nk}} - e^{-\sum_k l_{ik}\theta_{nk}}) + y_i, \end{aligned}$$

where S_{ij} is the set of pixels between the source and pixel j along projection i . If j' is the next pixel after pixel j along projection i , then

$$\begin{aligned} M_{ij} &= E(X_{ij} | Y_i = y_i, \boldsymbol{\theta}_n) \\ N_{ij} &= E(X_{ij'} | Y_i = y_i, \boldsymbol{\theta}_n). \end{aligned}$$

In view of expression (9.7), we find that

$$Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n) = \sum_i \sum_j \left[-N_{ij} l_{ij} \theta_j + (M_{ij} - N_{ij}) \ln(1 - e^{-l_{ij} \theta_j}) \right]$$

up to an irrelevant constant.

If we try to maximize $Q(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$ by setting its partial derivatives equal to 0, we get for pixel j the equation

$$-\sum_i N_{ij} l_{ij} + \sum_i \frac{(M_{ij} - N_{ij}) l_{ij}}{e^{l_{ij} \theta_j} - 1} = 0. \quad (9.8)$$

This is an intractable transcendental equation in the single variable θ_j , and the M step must be solved numerically, say by Newton's method. It is straightforward to check that the left-hand side of equation (9.8) is strictly decreasing in θ_j and has exactly one positive solution. Thus, the EM algorithm like the MM algorithm has the advantages of decoupling the parameters in the likelihood equations and of satisfying the natural boundary constraints $\theta_j \geq 0$. The MM algorithm is preferable to the EM algorithm because the MM algorithm involves far fewer exponentiations in defining its surrogate function.

9.7 Factor Analysis

In some instances, the missing data framework of the EM algorithm offers the easiest way to exploit convexity in deriving an MM algorithm. The complete data for a given problem is often fairly natural, and the difficulty

in deriving an EM algorithm shifts toward specifying the E step. Statisticians are particularly adept at calculating complicated conditional expectations connected with sampling distributions. We now illustrate these truths for estimation in factor analysis. Factor analysis explains the covariation among the components of a random vector by approximating the vector by a linear transformation of a small number of uncorrelated factors. Because factor analysis models usually involve normally distributed random vectors, Appendix A.2 reviews some basic facts about the multivariate normal distribution.

For the sake of notational convenience, we now extend the expectation and variance operators to random vectors. The expectation of a random vector $\mathbf{X} = (X_1, \dots, X_n)^*$ is defined componentwise by

$$\mathbf{E}(\mathbf{X}) = \begin{pmatrix} \mathbf{E}[X_1] \\ \vdots \\ \mathbf{E}[X_n] \end{pmatrix}.$$

Linearity carries over from the scalar case in the sense that

$$\begin{aligned} \mathbf{E}(\mathbf{X} + \mathbf{Y}) &= \mathbf{E}(\mathbf{X}) + \mathbf{E}(\mathbf{Y}) \\ \mathbf{E}(\mathbf{M}\mathbf{X}) &= \mathbf{M}\mathbf{E}(\mathbf{X}) \end{aligned}$$

for a compatible random vector \mathbf{Y} and a compatible matrix \mathbf{M} . The same componentwise conventions hold for the expectation of a random matrix and the variances and covariances of a random vector. Thus, we can express the variance matrix of a random vector \mathbf{X} as

$$\text{Var}(\mathbf{X}) = \mathbf{E}\{[\mathbf{X} - \mathbf{E}(\mathbf{X})][\mathbf{X} - \mathbf{E}(\mathbf{X})]^*\} = \mathbf{E}(\mathbf{X}\mathbf{X}^*) - \mathbf{E}(\mathbf{X})\mathbf{E}(\mathbf{X})^*.$$

These notational choices produce many other compact formulas. For instance, the random quadratic form $\mathbf{X}^*\mathbf{M}\mathbf{X}$ has expectation

$$\mathbf{E}(\mathbf{X}^*\mathbf{M}\mathbf{X}) = \text{tr}[\mathbf{M}\text{Var}(\mathbf{X})] + \mathbf{E}(\mathbf{X})^*\mathbf{M}\mathbf{E}(\mathbf{X}). \quad (9.9)$$

To verify this assertion, observe that

$$\begin{aligned} \mathbf{E}(\mathbf{X}^*\mathbf{M}\mathbf{X}) &= \mathbf{E}\left(\sum_i \sum_j X_i m_{ij} X_j\right) \\ &= \sum_i \sum_j m_{ij} \mathbf{E}(X_i X_j) \\ &= \sum_i \sum_j m_{ij} [\text{Cov}(X_i, X_j) + \mathbf{E}(X_i)\mathbf{E}(X_j)] \\ &= \text{tr}[\mathbf{M}\text{Var}(\mathbf{X})] + \mathbf{E}(\mathbf{X})^*\mathbf{M}\mathbf{E}(\mathbf{X}). \end{aligned}$$

The classical factor analysis model deals with l independent multivariate observations of the form

$$\mathbf{Y}_k = \boldsymbol{\mu} + \mathbf{F}\mathbf{X}_k + \mathbf{U}_k.$$

Here the $p \times q$ factor loading matrix \mathbf{F} transforms the unobserved factor score \mathbf{X}_k into the observed \mathbf{Y}_k . The random vector \mathbf{U}_k represents random measurement error. Typically, q is much smaller than p . The random vectors \mathbf{X}_k and \mathbf{U}_k are independent and normally distributed with means and variances

$$\begin{aligned} \mathbf{E}(\mathbf{X}_k) &= \mathbf{0}, & \text{Var}(\mathbf{X}_k) &= \mathbf{I} \\ \mathbf{E}(\mathbf{U}_k) &= \mathbf{0}, & \text{Var}(\mathbf{U}_k) &= \mathbf{D}, \end{aligned}$$

where \mathbf{I} is the $q \times q$ identity matrix and \mathbf{D} is a $p \times p$ diagonal matrix with i th diagonal entry d_i . The entries of the mean vector $\boldsymbol{\mu}$, the factor loading matrix \mathbf{F} , and the diagonal matrix \mathbf{D} constitute the parameters of the model. For a particular realization $\mathbf{y}_1, \dots, \mathbf{y}_l$ of the model, the maximum likelihood estimation of $\boldsymbol{\mu}$ is simply the sample mean $\hat{\boldsymbol{\mu}} = \bar{\mathbf{y}}$. This fact is a consequence of the reasoning given in Example 6.5.7. Therefore, we replace each \mathbf{y}_k by $\mathbf{y}_k - \bar{\mathbf{y}}$, assume $\boldsymbol{\mu} = \mathbf{0}$, and focus on estimating \mathbf{F} and \mathbf{D} .

The random vector $(\mathbf{X}_k^*, \mathbf{Y}_k^*)^*$ is the obvious choice of the complete data for case k . If $f(\mathbf{x}_k)$ is the density of \mathbf{X}_k and $g(\mathbf{y}_k | \mathbf{x}_k)$ is the conditional density of \mathbf{Y}_k given $\mathbf{X}_k = \mathbf{x}_k$, then the complete data loglikelihood can be expressed as

$$\begin{aligned} & \sum_{k=1}^l \ln f(\mathbf{x}_k) + \sum_{k=1}^l \ln g(\mathbf{y}_k | \mathbf{x}_k) \\ &= -\frac{l}{2} \ln \det \mathbf{I} - \frac{1}{2} \sum_{k=1}^l \mathbf{x}_k^* \mathbf{x}_k - \frac{l}{2} \ln \det \mathbf{D} \\ & \quad - \frac{1}{2} \sum_{k=1}^l (\mathbf{y}_k - \mathbf{F} \mathbf{x}_k)^* \mathbf{D}^{-1} (\mathbf{y}_k - \mathbf{F} \mathbf{x}_k). \end{aligned} \quad (9.10)$$

We can simplify this by noting that $\ln \det \mathbf{I} = 0$ and $\ln \det \mathbf{D} = \sum_{i=1}^p \ln d_i$.

The key to performing the E step is to note that $(\mathbf{X}_k^*, \mathbf{Y}_k^*)^*$ follows a multivariate normal distribution with variance matrix

$$\text{Var} \begin{pmatrix} \mathbf{X}_k \\ \mathbf{Y}_k \end{pmatrix} = \begin{pmatrix} \mathbf{I} & \mathbf{F}^* \\ \mathbf{F} & \mathbf{F} \mathbf{F}^* + \mathbf{D} \end{pmatrix}.$$

Equation (A.1) of Appendix A.2 then permits us to calculate the conditional expectation

$$\begin{aligned} \mathbf{v}_k &= \mathbf{E}(\mathbf{X}_k | \mathbf{Y}_k = \mathbf{y}_k, \mathbf{F}_n, \mathbf{D}_n) \\ &= \mathbf{F}_n^* (\mathbf{F}_n \mathbf{F}_n^* + \mathbf{D}_n)^{-1} \mathbf{y}_k \end{aligned}$$

and the conditional variance

$$\begin{aligned} \mathbf{A}_k &= \text{Var}(\mathbf{X}_k | \mathbf{Y}_k = \mathbf{y}_k, \mathbf{F}_n, \mathbf{D}_n) \\ &= \mathbf{I} - \mathbf{F}_n^* (\mathbf{F}_n \mathbf{F}_n^* + \mathbf{D}_n)^{-1} \mathbf{F}_n, \end{aligned}$$

given the observed data and the current values of the matrices \mathbf{F} and \mathbf{D} . Combining these results with equation (9.9) yields

$$\begin{aligned} & \mathbb{E}[(\mathbf{Y}_k - \mathbf{F}\mathbf{X}_k)^* \mathbf{D}^{-1}(\mathbf{Y}_k - \mathbf{F}\mathbf{X}_k) \mid \mathbf{Y}_k = \mathbf{y}_k] \\ &= \text{tr}(\mathbf{D}^{-1} \mathbf{F} \mathbf{A}_k \mathbf{F}^*) + (\mathbf{y}_k - \mathbf{F}\mathbf{v}_k)^* \mathbf{D}^{-1}(\mathbf{y}_k - \mathbf{F}\mathbf{v}_k) \\ &= \text{tr}\{\mathbf{D}^{-1}[\mathbf{F} \mathbf{A}_k \mathbf{F}^* + (\mathbf{y}_k - \mathbf{F}\mathbf{v}_k)(\mathbf{y}_k - \mathbf{F}\mathbf{v}_k)^*]\}. \end{aligned}$$

If we define

$$\mathbf{\Lambda} = \sum_{k=1}^l [\mathbf{A}_k + \mathbf{v}_k \mathbf{v}_k^*], \quad \mathbf{\Gamma} = \sum_{k=1}^l \mathbf{v}_k \mathbf{y}_k^*, \quad \mathbf{\Omega} = \sum_{k=1}^l \mathbf{y}_k \mathbf{y}_k^*$$

and take conditional expectations in equation (9.10), then we can write the surrogate function of the E step as

$$\begin{aligned} & Q(\mathbf{F}, \mathbf{D} \mid \mathbf{F}_n, \mathbf{D}_n) \\ &= -\frac{l}{2} \sum_{i=1}^p \ln d_i - \frac{1}{2} \text{tr}[\mathbf{D}^{-1}(\mathbf{F} \mathbf{\Lambda} \mathbf{F}^* - \mathbf{F} \mathbf{\Gamma} - \mathbf{\Gamma}^* \mathbf{F}^* + \mathbf{\Omega})], \end{aligned}$$

omitting the additive constant

$$-\frac{1}{2} \sum_{k=1}^l \mathbb{E}(\mathbf{X}_k^* \mathbf{X}_k \mid \mathbf{Y}_k = \mathbf{y}_k, \mathbf{F}_n, \mathbf{D}_n),$$

which depends on neither \mathbf{F} nor \mathbf{D} .

To perform the M step, we first maximize $Q(\mathbf{F}, \mathbf{D} \mid \mathbf{F}_n, \mathbf{D}_n)$ with respect to \mathbf{F} , holding \mathbf{D} fixed. We can do so by permuting factors and completing the square in the trace

$$\begin{aligned} & \text{tr}[\mathbf{D}^{-1}(\mathbf{F} \mathbf{\Lambda} \mathbf{F}^* - \mathbf{F} \mathbf{\Gamma} - \mathbf{\Gamma}^* \mathbf{F}^* + \mathbf{\Omega})] \\ &= \text{tr}[\mathbf{D}^{-1}(\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1}) \mathbf{\Lambda} (\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1})^*] + \text{tr}[\mathbf{D}^{-1}(\mathbf{\Omega} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1} \mathbf{\Gamma})] \\ &= \text{tr}[\mathbf{D}^{-\frac{1}{2}}(\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1}) \mathbf{\Lambda} (\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1})^* \mathbf{D}^{-\frac{1}{2}}] + \text{tr}[\mathbf{D}^{-1}(\mathbf{\Omega} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1} \mathbf{\Gamma})]. \end{aligned}$$

This calculation depends on the existence of the inverse matrix $\mathbf{\Lambda}^{-1}$. Now $\mathbf{\Lambda}$ is certainly positive definite if \mathbf{A}_k is positive definite, and Problem 22 asserts that \mathbf{A}_k is positive definite. It follows that $\mathbf{\Lambda}^{-1}$ not only exists but is positive definite as well. Furthermore, the matrix

$$\mathbf{D}^{-\frac{1}{2}}(\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1}) \mathbf{\Lambda} (\mathbf{F} - \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1})^* \mathbf{D}^{-\frac{1}{2}}$$

is positive semidefinite and has a nonnegative trace. Hence, the maximum value of the surrogate function $Q(\mathbf{F}, \mathbf{D} \mid \mathbf{F}_n, \mathbf{D}_n)$ with respect to \mathbf{F} is attained at the point $\mathbf{F} = \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1}$, regardless of the value of \mathbf{D} . In other words, the EM update of \mathbf{F} is $\mathbf{F}_{n+1} = \mathbf{\Gamma}^* \mathbf{\Lambda}^{-1}$. It should be stressed that

$\mathbf{\Gamma}$ and $\mathbf{\Lambda}$ implicitly depend on the previous values \mathbf{F}_n and \mathbf{D}_n . Once \mathbf{F}_{n+1} is determined, the equation

$$\begin{aligned} 0 &= \frac{\partial}{\partial d_i} Q(\mathbf{F}, \mathbf{D} \mid \mathbf{F}_n, \mathbf{D}_n) \\ &= -\frac{l}{2d_i} + \frac{1}{2d_i^2} (\mathbf{F}\mathbf{\Lambda}\mathbf{F}^* - \mathbf{F}\mathbf{\Gamma} - \mathbf{\Gamma}^*\mathbf{F}^* + \mathbf{\Omega})_{ii} \end{aligned}$$

provides the update

$$d_{n+1,i} = \frac{1}{l} (\mathbf{F}_{n+1}\mathbf{\Lambda}\mathbf{F}_{n+1}^* - \mathbf{F}_{n+1}\mathbf{\Gamma} - \mathbf{\Gamma}^*\mathbf{F}_{n+1}^* + \mathbf{\Omega})_{ii}.$$

One of the frustrating features of factor analysis is that the factor loading matrix \mathbf{F} is not uniquely determined. To understand the source of the ambiguity, consider replacing \mathbf{F} by $\mathbf{F}\mathbf{O}$, where \mathbf{O} is a $q \times q$ orthogonal matrix. The distribution of each random vector \mathbf{Y}_k is normal with mean $\boldsymbol{\mu}$ and variance matrix $\mathbf{F}\mathbf{F}^* + \mathbf{D}$. If we substitute $\mathbf{F}\mathbf{O}$ for \mathbf{F} , then the variance $\mathbf{F}\mathbf{O}\mathbf{O}^*\mathbf{F}^* + \mathbf{D} = \mathbf{F}\mathbf{F}^* + \mathbf{D}$ remains the same. Another problem in factor analysis is the existence of more than one local maximum. Which one of these the EM algorithm converges to depends on its starting value [76]. For a suggestion of how to improve the chances of converging to the dominant mode, see the article [281].

9.8 Hidden Markov Chains

A hidden Markov chain incorporates both observed data and missing data. The missing data are the sequence of states visited by the chain; the observed data provide partial information about this sequence of states. Denote the sequence of visited states by Z_1, \dots, Z_n and the observation taken at epoch i when the chain is in state Z_i by $Y_i = y_i$. Baum's algorithms [8, 71] recursively compute the likelihood of the observed data

$$P = \Pr(Y_1 = y_1, \dots, Y_n = y_n) \quad (9.11)$$

without actually enumerating all possible realizations Z_1, \dots, Z_n . Baum's algorithms can be adapted to perform an EM search. The references [78, 165, 216] discuss several concrete examples of hidden Markov chains.

The likelihood (9.11) is constructed from three ingredients: (a) the initial distribution $\boldsymbol{\pi}$ at the first epoch of the chain, (b) the epoch-dependent transition probabilities $p_{ijk} = \Pr(Z_{i+1} = k \mid Z_i = j)$, and (c) the conditional densities $\phi_i(y_i \mid j) = \Pr(Y_i = y_i \mid Z_i = j)$. The dependence of the transition probability p_{ijk} on i allows the chain to be inhomogeneous over time and promotes greater flexibility in modeling. Implicit in the definition of $\phi_i(y_i \mid j)$ are the assumptions that Y_1, \dots, Y_n are independent given

Z_1, \dots, Z_n and that Y_i depends only on Z_i . For simplicity, we will assume that the Y_i are discretely distributed.

Baum's forward algorithm is based on recursively evaluating the joint probabilities

$$\alpha_i(j) = \Pr(Y_1 = y_1, \dots, Y_{i-1} = y_{i-1}, Z_i = j).$$

At the first epoch, $\alpha_1(j) = \pi_j$ by definition. The obvious update to $\alpha_i(j)$ is

$$\alpha_{i+1}(k) = \sum_j \alpha_i(j) \phi_i(y_i | j) p_{ijk}. \tag{9.12}$$

The likelihood (9.11) can be recovered by computing the sum

$$P = \sum_j \alpha_n(j) \phi_n(y_n | j)$$

at the final epoch n .

In Baum's backward algorithm, we recursively evaluate the conditional probabilities

$$\beta_i(k) = \Pr(Y_{i+1} = y_{i+1}, \dots, Y_n = y_n | Z_i = k),$$

starting by convention at $\beta_n(k) = 1$ for all k . The required update is clearly

$$\beta_i(j) = \sum_k p_{ijk} \phi_{i+1}(y_{i+1} | k) \beta_{i+1}(k). \tag{9.13}$$

In this instance, the likelihood is recovered at the first epoch by forming the sum $P = \sum_j \pi_j \phi_1(y_1 | j) \beta_1(j)$.

Baum's algorithms also interdigitate beautifully with the E step of the EM algorithm. It is natural to summarize the missing data by a collection of indicator random variables X_{ij} . If the chain occupies state j at epoch i , then we take $X_{ij} = 1$. Otherwise, we take $X_{ij} = 0$. In this notation, the complete data loglikelihood can be written as

$$\begin{aligned} L_{\text{com}}(\boldsymbol{\theta}) &= \sum_j X_{1j} \ln \pi_j + \sum_{i=1}^n \sum_j X_{ij} \ln \phi_i(Y_i | j) \\ &\quad + \sum_{i=1}^{n-1} \sum_j \sum_k X_{ij} X_{i+1,k} \ln p_{ijk}. \end{aligned}$$

Execution of the E step amounts to calculation of the conditional expectations

$$E(X_{ij} X_{i+1,k} | \mathbf{Y}, \boldsymbol{\theta}_m) = \frac{\alpha_i(j) \phi_i(y_i | j) p_{ijk} \phi_{i+1}(y_{i+1} | k) \beta_{i+1}(k)}{P} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_m}$$

$$E(X_{ij} | \mathbf{Y}, \boldsymbol{\theta}_m) = \frac{\alpha_i(j) \phi_i(y_i | j) \beta_i(j)}{P} \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_m},$$

where $\mathbf{Y} = \mathbf{y}$ is the observed data, P is the likelihood of the observed data, and $\boldsymbol{\theta}_m$ is the current parameter vector.

The M step may or may not be exactly solvable. If it is not, then one can always revert to the MM gradient algorithm discussed in Sect. 10.4. In the case of hidden multinomial trials, it is possible to carry out the M step analytically. Hidden multinomial trials may govern (a) the choice of the initial state j , (b) the choice of an observed outcome Y_i at the i th epoch given the hidden state j of the chain at that epoch, or (c) the choice of the next state k given the current state j in a time-homogeneous chain. In the first case, the multinomial parameters are the π_j ; in the last case, they are the common transition probabilities p_{jk} .

As a concrete example, consider estimation of the initial distribution $\boldsymbol{\pi}$ at the first epoch of the chain. For estimation to be accurate, there must be multiple independent runs of the chain. Let the superscript r index the various runs. The surrogate function delivered by the E step equals

$$Q(\boldsymbol{\pi} \mid \boldsymbol{\pi}_m) = \sum_r \sum_j \mathbb{E}(X_{1j}^r \mid \mathbf{Y}^r = \mathbf{y}^r, \boldsymbol{\pi}_m) \ln \pi_j$$

up to an additive constant. Maximizing $Q(\boldsymbol{\pi} \mid \boldsymbol{\pi}_m)$ subject to the constraints $\sum_j \pi_j = 1$ and $\pi_j \geq 0$ for all j is done as in Example 1.4.2. The resulting EM updates

$$\pi_{m+1,j} = \frac{\sum_r \mathbb{E}(X_{1j}^r \mid \mathbf{Y}^r = \mathbf{y}^r, \boldsymbol{\pi}_m)}{R}$$

for R runs can be interpreted as multinomial proportions with fractional category counts. Problem 24 asks the reader to derive the EM algorithm for estimating time homogeneous transition probabilities. Problem 25 covers estimation of the parameters of the conditional densities $\phi_i(y_i \mid j)$ for some common densities.

9.9 Problems

1. Code and test any of the algorithms discussed in the text or problems of this chapter.
2. The entropy of a probability density $p(\mathbf{x})$ on \mathbb{R}^n is defined by

$$- \int p(\mathbf{x}) \ln p(\mathbf{x}) d\mathbf{x}. \quad (9.14)$$

Among all densities with a fixed mean vector $\boldsymbol{\mu} = \int \mathbf{x} p(\mathbf{x}) d\mathbf{x}$ and variance matrix $\boldsymbol{\Omega} = \int (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^* p(\mathbf{x}) d\mathbf{x}$, prove that the multivariate normal has maximum entropy. (Hints: Apply Proposition 9.2.1 and formula (9.9).)

3. In statistical mechanics, entropy is employed to characterize the stationary distribution of many independently behaving particles. Let $p(\mathbf{x})$ be the probability density that a particle is found at position \mathbf{x} in phase space \mathbb{R}^n , and suppose that each position \mathbf{x} is assigned an energy $u(\mathbf{x})$. If the average energy $U = \int u(\mathbf{x})p(\mathbf{x})d\mathbf{x}$ per particle is fixed, then Nature chooses $p(\mathbf{x})$ to maximize entropy as defined in equation (9.14). Show that if constants α and β exist satisfying

$$\int \alpha e^{\beta u(\mathbf{x})} d\mathbf{x} = 1 \quad \text{and} \quad \int u(\mathbf{x})\alpha e^{\beta u(\mathbf{x})} d\mathbf{x} = U,$$

then $p(\mathbf{x}) = \alpha e^{\beta u(\mathbf{x})}$ does indeed maximize entropy subject to the average energy constraint. The density $p(\mathbf{x})$ is the celebrated Maxwell-Boltzmann density.

4. Show that the normal, Poisson, binomial, negative binomial, gamma, beta, and multinomial families are exponential by writing their densities in the form (9.1). What are the corresponding measure and sufficient statistic in each case?
5. In the EM algorithm [65], suppose that the complete data X possesses a regular exponential density

$$f(x | \boldsymbol{\theta}) = g(x)e^{\beta(\boldsymbol{\theta})+h(x)*\boldsymbol{\theta}}$$

relative to some measure ν . Prove that the unconditional mean of the sufficient statistic $h(X)$ is given by the negative gradient $-\nabla\beta(\boldsymbol{\theta})$ and that the EM update is characterized by the condition

$$E[h(X) | Y, \theta_n] = -\nabla\beta(\boldsymbol{\theta}_{n+1}).$$

6. Suppose the phenotypic counts in the ABO allele frequency estimation example satisfy $n_A + n_{AB} > 0$, $n_B + n_{AB} > 0$, and $n_O > 0$. Show that the loglikelihood is strictly concave and possesses a single global maximum on the interior of the feasible region.
7. In a genetic linkage experiment, 197 animals are randomly assigned to four categories according to the multinomial distribution with cell probabilities $\pi_1 = \frac{1}{2} + \frac{\theta}{4}$, $\pi_2 = \frac{1-\theta}{4}$, $\pi_3 = \frac{1-\theta}{4}$ and $\pi_4 = \frac{\theta}{4}$. If the corresponding observations are

$$y = (y_1, y_2, y_3, y_4)^* = (125, 18, 20, 34)^*,$$

then devise an EM algorithm and use it to estimate $\hat{\theta} = .6268$ [218]. (Hint: Split the first category into two so that there are five categories for the complete data.)

8. Derive the EM algorithm solving Problem 7 as an MM algorithm. No mention of missing data is necessary.
9. Consider the data from *The London Times* [259] during the years 1910–1912 given in Table 9.1. The two columns labeled “Deaths i ” refer to the number of deaths to women 80 years and older reported by day. The columns labeled “Frequency n_i ” refer to the number of days with i deaths. A Poisson distribution gives a poor fit to these data, possibly because of different patterns of deaths in winter and summer. A mixture of two Poissons provides a much better fit. Under the Poisson admixture model, the likelihood of the observed data is

$$\prod_{i=0}^9 \left[\alpha e^{-\mu_1} \frac{\mu_1^i}{i!} + (1 - \alpha) e^{-\mu_2} \frac{\mu_2^i}{i!} \right]^{n_i},$$

where α is the admixture parameter and μ_1 and μ_2 are the means of the two Poisson distributions.

TABLE 9.1. Death notices from *The London Times*

Deaths i	Frequency n_i	Deaths i	Frequency n_i
0	162	5	61
1	267	6	27
2	271	7	8
3	185	8	3
4	111	9	1

Formulate an EM algorithm for this model. Let $\boldsymbol{\theta} = (\alpha, \mu_1, \mu_2)^*$ and

$$z_i(\boldsymbol{\theta}) = \frac{\alpha e^{-\mu_1} \mu_1^i}{\alpha e^{-\mu_1} \mu_1^i + (1 - \alpha) e^{-\mu_2} \mu_2^i}$$

be the posterior probability that a day with i deaths belongs to Poisson population 1. Show that the EM algorithm is given by

$$\begin{aligned} \alpha_{m+1} &= \frac{\sum_i n_i z_i(\boldsymbol{\theta}_m)}{\sum_i n_i} \\ \mu_{m+1,1} &= \frac{\sum_i n_i i z_i(\boldsymbol{\theta}_m)}{\sum_i n_i z_i(\boldsymbol{\theta}_m)} \\ \mu_{m+1,2} &= \frac{\sum_i n_i i [1 - z_i(\boldsymbol{\theta}_m)]}{\sum_i n_i [1 - z_i(\boldsymbol{\theta}_m)]}. \end{aligned}$$

From the initial estimates $\alpha_0 = 0.3$, $\mu_{01} = 1$, and $\mu_{02} = 2.5$, compute via the EM algorithm the maximum likelihood estimates $\hat{\alpha} = 0.3599$, $\hat{\mu}_1 = 1.2561$, and $\hat{\mu}_2 = 2.6634$. Note how slowly the EM algorithm converges in this example.

10. Derive the least squares algorithm (8.19) as an EM algorithm [112]. (Hint: Decompose y_i as the sum $\sum_{j=1}^q y_{ij}$ of realizations from independent normal deviates with means $x_{ij}\theta_j$ and variances $1/q$.)
11. Let x_1, \dots, x_m be an i.i.d. sample from a normal density with mean μ and variance σ^2 . Suppose for each x_i we observe $y_i = |x_i|$ rather than x_i . Formulate an EM algorithm for estimating μ and σ^2 , and show that its updates are

$$\begin{aligned}\mu_{n+1} &= \frac{1}{m} \sum_{i=1}^m (w_{ni1}y_i - w_{ni2}y_i) \\ \sigma_{n+1}^2 &= \frac{1}{m} \sum_{i=1}^m [w_{ni1}(y_i - \mu_{n+1})^2 + w_{ni2}(-y_i - \mu_{n+1})^2]\end{aligned}$$

with weights

$$\begin{aligned}w_{ni1} &= \frac{f(y_i | \boldsymbol{\theta}_n)}{f(y_i | \boldsymbol{\theta}_n) + f(-y_i | \boldsymbol{\theta}_n)} \\ w_{ni2} &= \frac{f(-y_i | \boldsymbol{\theta}_n)}{f(y_i | \boldsymbol{\theta}_n) + f(-y_i | \boldsymbol{\theta}_n)},\end{aligned}$$

where $f(x | \boldsymbol{\theta})$ is the normal density with $\boldsymbol{\theta} = (\mu, \sigma^2)^*$. Demonstrate that the modes of the likelihood of the observed data come in symmetric pairs differing only in the sign of μ . This fact does not prevent accurate estimation of $|\mu|$ and σ^2 .

12. Consider an i.i.d. sample drawn from a bivariate normal distribution with mean vector $\boldsymbol{\mu} = (\mu_1, \mu_2)^*$ and variance matrix

$$\boldsymbol{\Omega} = \begin{pmatrix} \sigma_1^2 & \sigma_{12} \\ \sigma_{12} & \sigma_2^2 \end{pmatrix}.$$

Suppose through some random accident that the first p observations are missing their first component, the next q observations are missing their second component, and the last r observations are complete. Design an EM algorithm to estimate the five mean and variance parameters, taking as complete data the original data before the accidental loss.

13. The standard linear regression model can be written in matrix notation as $\mathbf{X} = \mathbf{A}\boldsymbol{\beta} + \mathbf{U}$. Here \mathbf{X} is the $r \times 1$ vector of responses, \mathbf{A} is the $r \times s$ design matrix, $\boldsymbol{\beta}$ is the $s \times 1$ vector of regression coefficients, and \mathbf{U} is the $r \times 1$ normally distributed error vector with mean $\mathbf{0}$ and variance $\sigma^2 \mathbf{I}$. The responses are right censored if for each i there is a constant c_i such that only $Y_i = \min\{c_i, X_i\}$ is observed.

The EM algorithm offers a vehicle for estimating the parameter vector $\boldsymbol{\theta} = (\boldsymbol{\beta}^*, \sigma^2)^*$ in the presence of right censoring [65, 251]. Show that

$$\begin{aligned}\boldsymbol{\beta}_{n+1} &= (\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* \mathbf{E}(\mathbf{X} \mid \mathbf{Y}, \boldsymbol{\theta}_n) \\ \sigma_{n+1}^2 &= \frac{1}{r} \mathbf{E}[(\mathbf{X} - \mathbf{A}\boldsymbol{\beta}_{n+1})^*(\mathbf{X} - \mathbf{A}\boldsymbol{\beta}_{n+1}) \mid \mathbf{Y}, \boldsymbol{\theta}_n].\end{aligned}$$

To compute the conditional expectations appearing in these formulas, let \mathbf{a}_i be the i th row of \mathbf{A} and define

$$H(v) = \frac{\frac{1}{\sqrt{2\pi}} e^{-\frac{v^2}{2}}}{\frac{1}{\sqrt{2\pi}} \int_v^\infty e^{-\frac{w^2}{2}} dw}.$$

For a censored observation $y_i = c_i < \infty$, prove that

$$\begin{aligned}\mathbf{E}(X_i \mid Y_i = c_i, \boldsymbol{\theta}_n) &= \mathbf{a}_i \boldsymbol{\beta}_n + \sigma_n H\left(\frac{c_i - \mathbf{a}_i \boldsymbol{\beta}_n}{\sigma_n}\right) \\ \mathbf{E}(X_i^2 \mid Y_i = c_i, \boldsymbol{\theta}_n) &= (\mathbf{a}_i \boldsymbol{\beta}_n)^2 + \sigma_n^2 \\ &\quad + \sigma_n (c_i + \mathbf{a}_i \boldsymbol{\beta}_n) H\left(\frac{c_i - \mathbf{a}_i \boldsymbol{\beta}_n}{\sigma_n}\right).\end{aligned}$$

Use these formulas to complete the specification of the EM algorithm.

14. In the transmission tomography model it is possible to approximate the solution of equation (9.8) to good accuracy in certain situations. Verify the expansion

$$\frac{1}{e^s - 1} = \frac{1}{s} - \frac{1}{2} + \frac{s}{12} + O(s^2).$$

Using the approximation $1/(e^s - 1) \approx 1/s - 1/2$ for $s = l_{ij}\theta_j$, show that

$$\theta_{n+1,j} = \frac{\sum_i (M_{ij} - N_{ij})}{\frac{1}{2} \sum_i (M_{ij} + N_{ij}) l_{ij}}$$

results. Can you motivate this result heuristically?

15. Suppose that the complete data in the EM algorithm involve N binomial trials with success probability θ per trial. Here N can be random or fixed. If M trials result in success, then the complete data likelihood can be written as $\theta^M (1 - \theta)^{N-M} c$, where c is an irrelevant constant. The E step of the EM algorithm amounts to forming

$$Q(\theta \mid \boldsymbol{\theta}_n) = \mathbf{E}(M \mid \mathbf{Y}, \boldsymbol{\theta}_n) \ln \theta + \mathbf{E}(N - M \mid \mathbf{Y}, \boldsymbol{\theta}_n) \ln(1 - \theta) + \ln c.$$

The binomial trials are hidden because only a function \mathbf{Y} of them is directly observed. The brief derivation in Sect. 9.8 shows that the EM update amounts to

$$\theta_{n+1} = \frac{\mathbb{E}(M \mid \mathbf{Y}, \theta_n)}{\mathbb{E}(N \mid \mathbf{Y}, \theta_n)}.$$

Prove that this is equivalent to the update

$$\theta_{n+1} = \theta_n + \frac{\theta_n(1 - \theta_n)}{\mathbb{E}(N \mid \mathbf{Y}, \theta_n)} \frac{d}{d\theta} L(\theta_n),$$

where $L(\theta)$ is the loglikelihood of the observed data \mathbf{Y} [270]. (Hint: Apply identity (8.4) of Chap. 8.)

16. As an example of hidden binomial trials, consider a random sample of twin pairs. Let u of these pairs consist of male pairs, v consist of female pairs, and w consist of opposite sex pairs. A simple model to explain these data involves a random Bernoulli choice for each pair dictating whether it consists of identical or nonidentical twins. Suppose that identical twins occur with probability p and nonidentical twins with probability $1 - p$. Once the decision is made as to whether the twins are identical, then sexes are assigned to the twins. If the twins are identical, one assignment of sex is made. If the twins are nonidentical, then two independent assignments of sex are made. Suppose boys are chosen with probability q and girls with probability $1 - q$. Model these data as hidden binomial trials. Derive the EM algorithm for estimating p and q .
17. Chun Li has derived an EM update for hidden multinomial trials. Let N denote the number of hidden trials, θ_i the probability of outcome i of k possible outcomes, and $L(\boldsymbol{\theta})$ the loglikelihood of the observed data \mathbf{Y} . Derive the EM update

$$\theta_{n+1,i} = \theta_{ni} + \frac{\theta_{ni}}{\mathbb{E}(N \mid \mathbf{Y}, \boldsymbol{\theta}_n)} \left[\frac{\partial}{\partial \theta_i} L(\boldsymbol{\theta}_n) - \sum_{j=1}^k \theta_{nj} \frac{\partial}{\partial \theta_j} L(\boldsymbol{\theta}_n) \right]$$

following the reasoning of Problem 15.

18. In this problem you are asked to formulate models for hidden Poisson and exponential trials [270]. If the number of trials is N and the mean per trial is θ , then show that the EM update in the Poisson case is

$$\theta_{n+1} = \theta_n + \frac{\theta_n}{\mathbb{E}(N \mid \mathbf{Y}, \theta_n)} \frac{d}{d\theta} L(\theta_n)$$

and in the exponential case is

$$\theta_{n+1} = \theta_n + \frac{\theta_n^2}{\mathbb{E}(N \mid \mathbf{Y}, \theta_n)} \frac{d}{d\theta} L(\theta_n),$$

where $L(\theta)$ is the loglikelihood of the observed data \mathbf{Y} .

19. Suppose light bulbs have an exponential lifetime with mean θ . Two experiments are conducted. In the first, the lifetimes y_1, \dots, y_m of m independent bulbs are observed. In the second, p independent bulbs are observed to burn out before time t , and q independent bulbs are observed to burn out after time t . In other words, the lifetimes in the second experiment are both left and right censored. Construct an EM algorithm for finding the maximum likelihood estimate of θ [95].
20. In many discrete probability models, only data with positive counts are observed. Counts that are 0 are missing. Show that the likelihoods for the binomial, Poisson, and negative binomial models truncated at 0 amount to

$$\begin{aligned} L_1(p) &= \prod_i \frac{\binom{m_i}{x_i} p^{x_i} (1-p)^{m_i-x_i}}{1 - (1-p)^{m_i}} \\ L_2(\lambda) &= \prod_i \frac{\lambda^{x_i} e^{-\lambda}}{x_i! (1 - e^{-\lambda})} \\ L_3(p) &= \prod_i \frac{\binom{m_i+x_i-1}{x_i} (1-p)^{x_i} p^{m_i}}{1 - p^{m_i}}. \end{aligned}$$

For observation i of the binomial model, there are x_i successes out of m_i trials with success probability p per trial. For observation i of the negative binomial model, there are x_i failures before m_i required successes. For each model, devise an EM algorithm that fills in the missing observations by imputing a geometrically distributed number of truncated observations for every real observation. Show that the EM updates reduce to

$$\begin{aligned} p_{n+1} &= \frac{\sum_i x_i}{\sum_i \frac{m_i}{1 - (1-p_n)^{m_i}}} \\ \lambda_{n+1} &= \frac{\sum_i x_i}{\sum_i \frac{1}{1 - e^{-\lambda_n}}} \\ p_{n+1} &= \frac{\sum_i \frac{m_i}{1 - p_n^{m_i}}}{\sum_i (x_i + \frac{m_i}{1 - p_n^{m_i}})} \end{aligned}$$

for the three models.

21. Demonstrate that the EM updates of the previous problem can be derived as MM updates based on the minorization

$$-\ln(1-u) \geq -\ln(1-u_n) + \frac{u_n}{1-u_n} \ln \frac{u}{u_n}$$

for u and u_n in the interval $(0, 1)$. Prove this minorization first. (Hint: If you rearrange the minorization, then Proposition 9.2.1 applies.)

22. Suppose that Σ is a positive definite matrix. Prove that the matrix $\mathbf{I} - \mathbf{F}^*(\mathbf{F}\mathbf{F}^* + \Sigma)^{-1}\mathbf{F}$ is also positive definite. This result is used in the derivation of the EM algorithm in Sect. 9.7. (Hints: For readers familiar with the sweep operator of computational statistics, the simplest proof relies on applying Propositions 7.5.2 and 7.5.3 of the reference [166].)
23. A certain company asks consumers to rate movies on an integer scale from 1 to 5. Let M_i be the set of movies rated by person i . Denote the cardinality of M_i by $|M_i|$. Each rater does so in one of two modes that we will call “quirky” and “consensus”. In quirky mode, i has a private rating distribution $(q_{i1}, q_{i2}, q_{i3}, q_{i4}, q_{i5})$ that applies to every movie regardless of its intrinsic merit. In consensus mode, rater i rates movie j according to the distribution $(c_{j1}, c_{j2}, c_{j3}, c_{j4}, c_{j5})$ shared with all other raters in consensus mode. For every movie i rates, he or she makes a quirky decision with probability π_i and a consensus decision with probability $1 - \pi_i$. These decisions are made independently across raters and movies. If x_{ij} is the rating given to movie j by rater i , then prove that the likelihood of the data is

$$L = \prod_i \prod_{j \in M_i} [\pi_i q_{ix_{ij}} + (1 - \pi_i) c_{jx_{ij}}].$$

Once we estimate the parameters, we can rank the reliability of rater i by the estimate $\hat{\pi}_i$ and the popularity of movie j by its estimated average rating $\sum_k k \hat{c}_{jk}$.

If we choose the natural course of estimating the parameters by maximum likelihood, then it is possible to derive an EM or MM algorithm. From the right perspectives, these two algorithms coincide. Let n denote iteration number and w_{nij} the weight

$$w_{nij} = \frac{\pi_{ni} q_{nix_{ij}}}{\pi_{ni} q_{nix_{ij}} + (1 - \pi_{ni}) c_{njx_{ij}}}.$$

Derive either algorithm and show that it updates the parameters by

$$\begin{aligned} \pi_{n+1,i} &= \frac{1}{|M_i|} \sum_{j \in M_i} w_{nij} \\ q_{n+1,ix} &= \frac{\sum_{j \in M_i} \mathbf{1}_{\{x_{ij}=x\}} w_{nij}}{\sum_{j \in M_i} w_{nij}} \\ c_{n+1,jx} &= \frac{\sum_i \mathbf{1}_{\{x_{ij}=x\}} (1 - w_{nij})}{\sum_i (1 - w_{nij})}. \end{aligned}$$

These updates are easy to implement. Can you motivate them as ratios of expected counts?

24. In the hidden Markov chain model, suppose that the chain is time homogeneous with transition probabilities p_{jk} . Derive an EM algorithm for estimating the p_{jk} from one or more independent runs of the chain.
25. In the hidden Markov chain model, consider estimation of the parameters of the conditional densities $\phi_i(y_i | j)$ of the observed data y_1, \dots, y_n . When Y_i given $Z_i = j$ is Poisson distributed with mean μ_j , show that the EM algorithm updates μ_j by

$$\mu_{m+1,j} = \frac{\sum_{i=1}^n w_{mij} y_i}{\sum_{i=1}^n w_{mij}},$$

where the weight $w_{mij} = E(X_{ij} | Y, \boldsymbol{\mu}_m)$. Show that the same update applies when Y_i given $Z_i = i$ is exponentially distributed with mean μ_j or normally distributed with mean μ_j and common variance σ^2 . In the latter setting, demonstrate that the EM update of σ^2 is

$$\sigma_{m+1}^2 = \frac{\sum_{i=1}^n \sum_j w_{mij} (y_i - \mu_{m+1,j})^2}{\sum_{i=1}^n \sum_j w_{mij}}.$$

10

Newton's Method and Scoring

10.1 Introduction

Block relaxation and the MM algorithm are hardly the only methods of optimization. Newton's method is better known and more widely applied. Despite its defects, Newton's method is the gold standard for speed of convergence and forms the basis of most modern optimization algorithms in low dimensions. Its many variants seek to retain its fast convergence while taming its defects. The variants all revolve around the core idea of locally approximating the objective function by a strictly convex quadratic function. At each iteration the quadratic approximation is optimized. Safeguards are introduced to keep the iterates from veering toward irrelevant stationary points.

Statisticians are among the most avid consumers of optimization techniques. Statistics, like other scientific disciplines, has a special vocabulary. We will meet some of that vocabulary in this chapter as we discuss optimization methods important in computational statistics. Thus, we will take up Fisher's scoring algorithm and the Gauss-Newton method of nonlinear least squares. We have already encountered likelihood functions and the device of passing to loglikelihoods. In statistics, the gradient of the loglikelihood is called the score, and the negative of the second differential is called the observed information. One major advantage of maximizing the loglikelihood rather than the likelihood is that the loglikelihood, score, and observed information are all additive functions of independent observations.

10.2 Newton's Method and Root Finding

One of the virtues of Newton's method is that it is a root-finding technique as well as an optimization technique. Consider a function $f(\mathbf{x})$ mapping \mathbb{R}^n into \mathbb{R}^n , and suppose a root of $f(\mathbf{x}) = \mathbf{0}$ occurs at \mathbf{y} . If the slope matrix in the expansion

$$\begin{aligned}\mathbf{0} - f(\mathbf{x}) &= f(\mathbf{y}) - f(\mathbf{x}) \\ &= s(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})\end{aligned}$$

is invertible, then we can solve for \mathbf{y} as

$$\mathbf{y} = \mathbf{x} - s(\mathbf{y}, \mathbf{x})^{-1}f(\mathbf{x}).$$

In practice, \mathbf{y} is unknown and the slope $s(\mathbf{y}, \mathbf{x})$ is unavailable. However, if \mathbf{y} is close to \mathbf{x} , then $s(\mathbf{y}, \mathbf{x})$ should be close to $df(\mathbf{x})$. Thus, Newton's method iterates according to

$$\mathbf{x}_{m+1} = \mathbf{x}_m - df(\mathbf{x}_m)^{-1}f(\mathbf{x}_m). \quad (10.1)$$

Example 10.2.1 *Division without Dividing*

Forming the reciprocal of a number $a > 0$ is equivalent to solving for a root of the equation $f(x) = a - x^{-1}$. Newton's method (10.1) iterates according to

$$x_{m+1} = x_m - \frac{a - x_m^{-1}}{x_m^{-2}} = x_m(2 - ax_m),$$

which involves multiplication and subtraction but no division. If x_{m+1} is to be positive, then x_m must lie on the interval $(0, 2/a)$. If x_m does indeed reside there, then x_{m+1} will reside on the shorter interval $(0, 1/a)$ because the quadratic $x(2 - ax)$ attains its maximum of $1/a$ at $x = 1/a$. Furthermore, $x_{m+1} > x_m$ if and only if $2 - ax_m > 1$, and this latter inequality holds if and only if $x_m < 1/a$. Thus, starting on $(0, 1/a)$, the iterates x_m monotonically increase to their limit $1/a$. Starting on $[1/a, 2/a)$, the first iterate satisfies $x_1 \leq 1/a$, and subsequent iterates monotonically increase to $1/a$. Finally, starting outside $(0, 2/a)$ leads either to fixation at 0 or divergence to $-\infty$. ■

Example 10.2.2 *Extraction of n th Roots*

Newton's method can be used to extract square roots, cube roots, and so forth. Consider the function $f(x) = x^n - a$ for some integer $n > 1$ and $a > 0$. Newton's method amounts to the iteration scheme

$$x_{m+1} = x_m - \frac{x_m^n - a}{nx_m^{n-1}} = \frac{1}{n} \left[(n-1)x_m + \frac{a}{x_m^{n-1}} \right]. \quad (10.2)$$

TABLE 10.1. Newton's iterates for $x^4 - x^2$

m	x_m	x_m	x_m
0	-0.74710	-0.66710	-0.500000
1	-2.16581	1.01669	-0.125000
2	-1.68896	1.00066	-0.061491
3	-1.35646	1.00000	-0.030628
4	-1.14388	1.00000	-0.015300
5	-1.03477	1.00000	-0.007648
6	-1.00270	1.00000	-0.003823
7	-1.00002	1.00000	-0.001911
8	-1.00000	1.00000	-0.000955
9	-1.00000	1.00000	-0.000477
10	-1.00000	1.00000	-0.000238

This sequence converges to $\sqrt[3]{a}$ regardless of the starting point $x_0 > 0$. To demonstrate this fact, we first note that the right-hand side of equation (10.2) is the arithmetic mean of $n-1$ copies of the number x_m and a/x_m^{n-1} . Because the arithmetic mean exceeds the geometric mean $\sqrt[3]{a}$, it follows that $x_m \geq \sqrt[3]{a}$ for all $m \geq 1$. Given this inequality, we have $a/x_m^{n-1} \leq x_m$. Again viewing equation (10.2) as a weighted average of x_m and the ratio $a/x_m^{n-1} \leq x_m$, it follows that $x_{m+1} \leq x_m$ for all $m \geq 1$. Hence, the sequence x_1, x_2, \dots is bounded below and is monotonically decreasing. By continuity, its limit is $\sqrt[3]{a}$. ■

Example 10.2.3 Sensitivity to Initial Conditions

In contrast to the previous two well-behaved examples, finding a root of the polynomial $f(x) = x^4 - x^2$ is more problematic. These roots are clearly $-1, 0$, and 1 . We anticipate trouble when $f'(x) = 4x^3 - 2x = 0$ at the points $-1/\sqrt{2}, 0$, and $1/\sqrt{2}$. Consider initial points near $-1/\sqrt{2}$. Just to the left of $-1/\sqrt{2}$, Newton's method converges to -1 . For a narrow zone just to the right of $-1/\sqrt{2}$, it converges to 1 , and beyond this zone but to the left of $1/\sqrt{2}$, it converges to 0 . Table 10.1 gives three typical examples of this extreme sensitivity to initial conditions. The slower convergence to the middle root 0 is hardly surprising given that $f'(0) = 0$. It appears that the discrepancy $x_m - 0$ roughly halves at each iteration. Problem 2 clarifies this behavior. ■

Example 10.2.4 Secant Method

There are several ways of estimating the differential $df(\mathbf{x}_m)$ appearing in Newton's formula (10.1). In one dimension the secant method approximates $f'(x_m)$ by the slope $s(x_{m-1}, x_m)$ using the canonical slope function

$$s(y, x) = \frac{f(y) - f(x)}{y - x}.$$

This produces the secant update

$$x_{m+1} = x_m - \frac{f(x_m)(x_{m-1} - x_m)}{f(x_{m-1}) - f(x_m)}, \quad (10.3)$$

which is the prototype for the quasi-Newton updates treated in Chap. 11. While the secant method avoids computation of derivatives, it typically takes more iterations to converge than Newton's method. Safeguards must also be put into place to ensure its reliability. ■

Example 10.2.5 *Newton's Method of Matrix Inversion*

Newton's method for finding the reciprocal of a number can be generalized to compute the inverse of a matrix [138]. Consider the matrix-valued function $f(\mathbf{B}) = \mathbf{A} - \mathbf{B}^{-1}$ for some invertible $n \times n$ matrix \mathbf{A} . Example 5.4.5 provides the first-order differential approximation

$$f(\mathbf{A}^{-1}) - f(\mathbf{B}) = \mathbf{0} - (\mathbf{A} - \mathbf{B}^{-1}) \approx \mathbf{B}^{-1}(\mathbf{A}^{-1} - \mathbf{B})\mathbf{B}^{-1}.$$

Multiplying this equation on the left and right by \mathbf{B} and rearranging give

$$\mathbf{A}^{-1} \approx 2\mathbf{B} - \mathbf{B}\mathbf{A}\mathbf{B}$$

and hence Newton's scheme

$$\mathbf{B}_{m+1} = 2\mathbf{B}_m - \mathbf{B}_m\mathbf{A}\mathbf{B}_m.$$

Further rearrangement yields

$$\mathbf{A}^{-1} - \mathbf{B}_{m+1} = (\mathbf{A}^{-1} - \mathbf{B}_m)\mathbf{A}(\mathbf{A}^{-1} - \mathbf{B}_m),$$

which entails

$$\|\mathbf{A}^{-1} - \mathbf{B}_{m+1}\| \leq \|\mathbf{A}\| \cdot \|\mathbf{A}^{-1} - \mathbf{B}_m\|^2$$

for every matrix norm. It follows that the sequence \mathbf{B}_m converges at a quadratic rate to \mathbf{A}^{-1} if \mathbf{B}_0 is sufficiently close to \mathbf{A}^{-1} . ■

10.3 Newton's Method and Optimization

Suppose we want to minimize the real-valued function $f(\mathbf{x})$ defined on an open set $S \subset \mathbb{R}^n$. Assuming that $f(\mathbf{x})$ is twice differentiable, the expansion

$$f(\mathbf{y}) = f(\mathbf{x}) + df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + \frac{1}{2}(\mathbf{y} - \mathbf{x})^* s^2(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x})$$

suggests that we substitute $d^2 f(\mathbf{x})$ for the second slope $s^2(\mathbf{y}, \mathbf{x})$ and approximate $f(\mathbf{y})$ by the resulting quadratic. If we take this approximation seriously, then we can solve for its minimum point \mathbf{y} as

$$\mathbf{y} = \mathbf{x} - d^2 f(\mathbf{x})^{-1} \nabla f(\mathbf{x}).$$

In Newton's method we iterate according to

$$\mathbf{x}_{m+1} = \mathbf{x}_m - d^2 f(\mathbf{x}_m)^{-1} \nabla f(\mathbf{x}_m). \tag{10.4}$$

It should come as no surprise that algorithm (10.4) coincides with the earlier version of Newton's method seeking a root of $\nabla f(\mathbf{x})$. For this reason, any stationary point \mathbf{x} of $f(\mathbf{x})$ is a fixed point of algorithm (10.4).

Example 10.3.1 *Newton's Method for the Poisson Multigraph Model*

Newton's method can be applied to the Poisson multigraph model introduced in Sect. 8.11. The score vector has entries

$$\frac{\partial}{\partial p_i} L(\mathbf{p}) = \sum_{j \neq i} \left(\frac{x_{ij}}{p_i} - p_j \right),$$

and the observed information matrix has entries

$$-\frac{\partial^2}{\partial p_i \partial p_j} L(\mathbf{p}) = \begin{cases} 1 & i \neq j \\ \frac{1}{p_i^2} \sum_{k \neq i} x_{ik} & i = j. \end{cases}$$

For n nodes the matrix $-d^2 L(\mathbf{p})$ is $n \times n$, and inverting it seems out of the question when n is large. Fortunately, the Sherman-Morrison formula comes to the rescue. If we write $-d^2 L(\mathbf{p})$ as $\mathbf{D} + \mathbf{1}\mathbf{1}^*$ with \mathbf{D} diagonal, then the explicit inverse

$$(\mathbf{D} + \mathbf{1}\mathbf{1}^*)^{-1} = \mathbf{D}^{-1} - \frac{1}{1 + \mathbf{1}^* \mathbf{D}^{-1} \mathbf{1}} \mathbf{D}^{-1} \mathbf{1}\mathbf{1}^* \mathbf{D}^{-1}$$

is available. This makes Newton's method trivial to implement as long as one respects the bounds $p_i \geq 0$. More generally, it is always cheap to invert a low-rank perturbation of an explicitly invertible matrix. See Problem 10 of Chap. 11 for Woodbury's generalization of the Sherman-Morrison formula.

■

There are two potential problems with Newton's method. First, it may be expensive computationally to evaluate or invert $d^2 f(\mathbf{x})$. Second, far from the minimum, Newton's method is equally happy to head uphill or down. In other words, Newton's method is not a descent algorithm in the sense that $f(\mathbf{x}_{m+1}) < f(\mathbf{x}_m)$. This second defect can be remedied by modifying the Newton increment so that it is a partial step in a descent direction. A descent direction \mathbf{v} at the point \mathbf{x} satisfies the inequality $df(\mathbf{x})\mathbf{v} < 0$. The formula

$$\lim_{t \downarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t} = df(\mathbf{x})\mathbf{v}$$

for the forward directional derivative shows that $f(\mathbf{x} + t\mathbf{v}) < f(\mathbf{x})$ for $t > 0$ sufficiently small. The key to generating a descent direction is to define

$\mathbf{v} = -\mathbf{H}^{-1}\nabla f(\mathbf{x})$ using a positive definite matrix \mathbf{H} . Here we assume that \mathbf{x} is not a stationary point and recall the fact that the inverse of a positive definite matrix is positive definite.

The necessary modifications of Newton's method to achieve a descent algorithm are now clear. We simply replace $d^2f(\mathbf{x}_m)$ by a positive definite approximating matrix \mathbf{H}_m and take a sufficiently short step in the direction $\Delta\mathbf{x}_m = -\mathbf{H}_m^{-1}\nabla f(\mathbf{x}_m)$. If we believe that the proposed increment is reasonable, then we will be reluctant to shrink $\Delta\mathbf{x}_m$ too much. This suggests backtracking, the simplest form of which is step halving. In step halving, if the initial increment $\Delta\mathbf{x}_m$ does not produce a decrease in $f(\mathbf{x})$, then try $\Delta\mathbf{x}_m/2$. If $\Delta\mathbf{x}_m/2$ fails, then try $\Delta\mathbf{x}_m/4$, and so forth. We will meet more sophisticated backtracking schemes later. Note at this juncture we have said nothing about how well \mathbf{H}_m approximates $d^2f(\mathbf{x}_m)$. The quality of this approximation obviously affects the rate of convergence toward any local minimum.

If we minimize $f(\mathbf{x})$ subject to the linear equality constraints $\mathbf{V}\mathbf{x} = \mathbf{d}$, then minimization of the approximating quadratic can be accomplished as indicated in Example 5.2.6 of Chap. 5. Because

$$\mathbf{V}(\mathbf{x}_{m+1} - \mathbf{x}_m) = \mathbf{0},$$

the revised increment $\Delta\mathbf{x}_m = \mathbf{x}_{m+1} - \mathbf{x}_m$ is

$$\Delta\mathbf{x}_m = -[\mathbf{H}_m^{-1} - \mathbf{H}_m^{-1}\mathbf{V}^*(\mathbf{V}\mathbf{H}_m^{-1}\mathbf{V}^*)^{-1}\mathbf{V}\mathbf{H}_m^{-1}]\nabla f(\mathbf{x}_m). \quad (10.5)$$

This can be viewed as the projection of the unconstrained increment onto the null space of \mathbf{V} . Problem 12 shows that step halving also works for the projected increment.

10.4 MM Gradient Algorithm

Often it is impossible to solve the optimization step of the MM algorithm exactly. If $f(\mathbf{x})$ is the objective function and $g(\mathbf{x} \mid \mathbf{x}_m)$ minorizes or majorizes $f(\mathbf{x})$ at \mathbf{x}_m , then Newton's method can be applied to optimize $g(\mathbf{x} \mid \mathbf{x}_m)$. As we shall see later, one step of Newton's method preserves the overall rate of convergence of the MM algorithm. Thus, the MM gradient algorithm iterates according to

$$\begin{aligned} \mathbf{x}_{m+1} &= \mathbf{x}_m - d^2g(\mathbf{x}_m \mid \mathbf{x}_m)^{-1}\nabla g(\mathbf{x}_m \mid \mathbf{x}_m) \\ &= \mathbf{x}_m - d^2g(\mathbf{x}_m \mid \mathbf{x}_m)^{-1}\nabla f(\mathbf{x}_m). \end{aligned}$$

Here derivatives are taken with respect to the left argument of $g(\mathbf{x} \mid \mathbf{x}_m)$. Substitution of $\nabla f(\mathbf{x}_m)$ for $\nabla g(\mathbf{x}_m \mid \mathbf{x}_m)$ is justified by the comments in Sect. 8.2. In practice the surrogate function $g(\mathbf{x} \mid \mathbf{x}_m)$ is either convex or concave, and its second differential $d^2g(\mathbf{x}_m \mid \mathbf{x}_m)$ needs no adjustment to give a descent or ascent algorithm.

Example 10.4.1 *Newton’s Method in Transmission Tomography*

In the transmission tomography model of Chap. 8, the surrogate function $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ of equation (8.15) minorizes the loglikelihood $L(\boldsymbol{\theta})$ in the absence of a smoothing prior. Differentiating $g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ with respect to θ_j gives the transcendental equation

$$0 = \sum_i l_{ij} \left[d_i e^{-\mathbf{l}_i^* \boldsymbol{\theta}_m \theta_j / \theta_{mj}} - y_i \right].$$

One step of Newton’s method starting at $\theta_j = \theta_{mj}$ produces the next iterate

$$\begin{aligned} \theta_{m+1,j} &= \theta_{mj} + \frac{\theta_{mj} \sum_i l_{ij} (d_i e^{-\mathbf{l}_i^* \boldsymbol{\theta}_m} - y_i)}{\sum_i l_{ij} \mathbf{l}_i^* \boldsymbol{\theta}_m d_i e^{-\mathbf{l}_i^* \boldsymbol{\theta}_m}} \\ &= \theta_{mj} \frac{\sum_i l_{ij} [d_i e^{-\mathbf{l}_i^* \boldsymbol{\theta}_m} (1 + \mathbf{l}_i^* \boldsymbol{\theta}_m) - y_i]}{\sum_i l_{ij} \mathbf{l}_i^* \boldsymbol{\theta}_m d_i e^{-\mathbf{l}_i^* \boldsymbol{\theta}_m}}. \end{aligned}$$

This step typically increases $L(\boldsymbol{\theta})$. The comparable EM gradient update involves solving the transcendental equation (9.8). ■

Example 10.4.2 *Estimation with the Dirichlet Distribution*

As another example, consider parameter estimation for the Dirichlet distribution [154]. This distribution has probability density

$$\frac{\Gamma(\sum_{i=1}^n \theta_i)}{\prod_{i=1}^n \Gamma(\theta_i)} \prod_{i=1}^n y_i^{\theta_i - 1} \tag{10.6}$$

on the unit simplex $\{\mathbf{y} = (y_1, \dots, y_n)^* : y_1 > 0, \dots, y_n > 0, \sum_{i=1}^n y_i = 1\}$ endowed with the uniform measure. The Dirichlet distribution is used to represent random proportions. The beta distribution is the special case $n = 2$.

If $\mathbf{y}_1, \dots, \mathbf{y}_l$ are randomly sampled vectors from the Dirichlet distribution, then their loglikelihood is

$$L(\boldsymbol{\theta}) = l \ln \Gamma\left(\sum_{i=1}^n \theta_i\right) - l \sum_{i=1}^n \ln \Gamma(\theta_i) + \sum_{k=1}^l \sum_{i=1}^n (\theta_i - 1) \ln y_{ki}.$$

Except for the first term on the right, the parameters are separated. Fortunately as demonstrated in Example 6.3.11, the function $\ln \Gamma(t)$ is convex. Denoting its derivative by $\psi(t)$, we exploit the minorization

$$\ln \Gamma\left(\sum_{i=1}^n \theta_i\right) \geq \ln \Gamma\left(\sum_{i=1}^n \theta_{mi}\right) + \psi\left(\sum_{i=1}^n \theta_{mi}\right) \sum_{i=1}^n (\theta_i - \theta_{mi})$$

and create the surrogate function

$$g(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = l \ln \Gamma\left(\sum_{i=1}^n \theta_{mi}\right) + l\psi\left(\sum_{i=1}^n \theta_{mi}\right) \sum_{i=1}^n (\theta_i - \theta_{mi}) \\ - l \sum_{i=1}^n \ln \Gamma(\theta_i) + \sum_{k=1}^l \sum_{i=1}^n (\theta_i - 1) \ln y_{ki}.$$

Owing to the presence of the terms $\ln \Gamma(\theta_i)$, the maximization step is intractable. However, the MM gradient algorithm can be readily implemented because the parameters are now separated and the functions $\psi(t)$ and $\psi'(t)$ are easily computed as suggested in Problem 14. The whole process is carried out in the references [163, 199] on actual data. ■

10.5 Ad Hoc Approximations of $d^2 f(\boldsymbol{\theta})$

In minimization problems, we have emphasized the importance of approximating $d^2 f(\boldsymbol{\theta})$ by a positive definite matrix. Three key ideas drive the process of approximation. One is the recognition that outer product matrices are positive semidefinite. Another is a feel for when terms are small on average. Usually this involves comparison of random variables and their means. Finally, it is almost always advantageous to avoid the explicit calculation of complicated second derivatives.

For example, consider the problem of least squares estimation with nonlinear regression functions. Let us formulate the problem slightly more generally as one of minimizing the sum of squares

$$f(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta})]^2$$

involving a weight $w_i > 0$ and response y_i for each case i . Here y_i is a realization of a random variable Y_i with mean $\mu_i(\boldsymbol{\theta})$. In linear regression, $\mu_i(\boldsymbol{\theta}) = \sum_k x_{ik} \theta_k$. To implement Newton's method, we need

$$\nabla f(\boldsymbol{\theta}) = - \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta})] \nabla \mu_i(\boldsymbol{\theta}) \\ d^2 f(\boldsymbol{\theta}) = \sum_{i=1}^n w_i \nabla \mu_i(\boldsymbol{\theta}) d\mu_i(\boldsymbol{\theta}) - \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta})] d^2 \mu_i(\boldsymbol{\theta}). \quad (10.7)$$

In the Gauss-Newton algorithm, we approximate

$$d^2 f(\boldsymbol{\theta}) \approx \sum_{i=1}^n w_i \nabla \mu_i(\boldsymbol{\theta}) d\mu_i(\boldsymbol{\theta})$$

on the rationale that either the weighted residuals $w_i[y_i - \mu_i(\theta)]$ are small or the regression functions $\mu_i(\theta)$ are nearly linear. In both instances, the Gauss-Newton algorithm shares the fast convergence of Newton's method.

Maximum likelihood estimation with the Poisson distribution furnishes another example. Here the count data y_1, \dots, y_n have loglikelihood, score, and negative observed information

$$\begin{aligned} L(\theta) &= \sum_{i=1}^n [y_i \ln \lambda_i(\theta) - \lambda_i(\theta) - \ln y_i!] \\ \nabla L(\theta) &= \sum_{i=1}^n \left[\frac{y_i}{\lambda_i(\theta)} \nabla \lambda_i(\theta) - \nabla \lambda_i(\theta) \right] \\ d^2 L(\theta) &= \sum_{i=1}^n \left[-\frac{y_i}{\lambda_i(\theta)^2} \nabla \lambda_i(\theta) d\lambda_i(\theta) + \frac{y_i}{\lambda_i(\theta)} d^2 \lambda_i(\theta) - d^2 \lambda_i(\theta) \right], \end{aligned}$$

where $E(y_i) = \lambda_i(\theta)$. Given that the ratio $y_i/\lambda_i(\theta)$ has average value 1, the negative semidefinite approximations

$$\begin{aligned} d^2 L(\theta) &\approx -\sum_{i=1}^n \frac{y_i}{\lambda_i(\theta)^2} \nabla \lambda_i(\theta) d\lambda_i(\theta) \\ &\approx -\sum_{i=1}^n \frac{1}{\lambda_i(\theta)} \nabla \lambda_i(\theta) d\lambda_i(\theta) \end{aligned}$$

are reasonable. The second of these leads to the scoring algorithm discussed in the next section.

The exponential distribution offers a third illustration. Now the data have means $E(y_i) = 1/\lambda_i(\theta)$. The loglikelihood

$$L(\theta) = \sum_{i=1}^n [\ln \lambda_i(\theta) - y_i \lambda_i(\theta)]$$

yields the score and negative observed information

$$\begin{aligned} \nabla L(\theta) &= \sum_{i=1}^n \left[\frac{1}{\lambda_i(\theta)} \nabla \lambda_i(\theta) - y_i \nabla \lambda_i(\theta) \right] \\ d^2 L(\theta) &= \sum_{i=1}^n \left[-\frac{1}{\lambda_i(\theta)^2} \nabla \lambda_i(\theta) d\lambda_i(\theta) + \frac{1}{\lambda_i(\theta)} d^2 \lambda_i(\theta) - y_i d^2 \lambda_i(\theta) \right]. \end{aligned}$$

Replacing observations by their means suggests the approximation

$$d^2 L(\theta) \approx -\sum_{i=1}^n \frac{1}{\lambda_i(\theta)^2} \nabla \lambda_i(\theta) d\lambda_i(\theta)$$

made in the scoring algorithm. Table 10.2 summarizes the scoring algorithm with means $\mu_i(\boldsymbol{\theta})$ replacing intensities $\lambda_i(\boldsymbol{\theta})$.

Our final example involves maximum likelihood estimation with the multinomial distribution. The observations y_1, \dots, y_j are now cell counts over n independent trials. Cell i is assigned probability $p_i(\boldsymbol{\theta})$ and averages a total of $np_i(\boldsymbol{\theta})$ counts. The loglikelihood, score, and negative observed information amount to

$$\begin{aligned} L(\boldsymbol{\theta}) &= \sum_{i=1}^j y_i \ln p_i(\boldsymbol{\theta}) \\ \nabla L(\boldsymbol{\theta}) &= \sum_{i=1}^j \frac{y_i}{p_i(\boldsymbol{\theta})} \nabla p_i(\boldsymbol{\theta}) \\ d^2 L(\boldsymbol{\theta}) &= \sum_{i=1}^j \left[-\frac{y_i}{p_i(\boldsymbol{\theta})^2} \nabla p_i(\boldsymbol{\theta}) dp_i(\boldsymbol{\theta}) + \frac{y_i}{p_i(\boldsymbol{\theta})} d^2 p_i(\boldsymbol{\theta}) \right]. \end{aligned}$$

In light of the identity $E(y_i) = np_i(\boldsymbol{\theta})$, the approximation

$$\sum_{i=1}^j \frac{y_i}{p_i(\boldsymbol{\theta})} d^2 p_i(\boldsymbol{\theta}) \approx n \sum_{i=1}^j d^2 p_i(\boldsymbol{\theta}) = nd^2 \mathbf{1} = \mathbf{0}$$

is reasonable. This suggests the further negative semidefinite approximations

$$\begin{aligned} d^2 L(\boldsymbol{\theta}) &\approx -\sum_{i=1}^j \frac{y_i}{p_i(\boldsymbol{\theta})^2} \nabla p_i(\boldsymbol{\theta}) dp_i(\boldsymbol{\theta}) \\ &\approx -n \sum_{i=1}^j \frac{1}{p_i(\boldsymbol{\theta})} \nabla p_i(\boldsymbol{\theta}) dp_i(\boldsymbol{\theta}), \end{aligned}$$

the second of which coincides with the scoring algorithm.

10.6 Scoring and Exponential Families

As we have just witnessed, one can approximate the observed information in a variety of ways. The method of steepest ascent replaces the observed information by the identity matrix \mathbf{I} . The usually more efficient scoring algorithm replaces the observed information by the expected information $J(\boldsymbol{\theta}) = E[-d^2 L(\boldsymbol{\theta})]$, where $L(\boldsymbol{\theta})$ is the loglikelihood. The alternative representation $J(\boldsymbol{\theta}) = \text{Var}[\nabla L(\boldsymbol{\theta})]$ of $J(\boldsymbol{\theta})$ as a variance matrix shows that it is positive semidefinite and hence a good replacement for $-d^2 L(\boldsymbol{\theta})$ in Newton's method. An extra dividend of scoring is that the inverse matrix $J(\hat{\boldsymbol{\theta}})^{-1}$ immediately supplies the asymptotic variances and covariances of

the maximum likelihood estimate $\hat{\boldsymbol{\theta}}$ [218]. Scoring shares this benefit with Newton's method since the observed information is under natural assumptions asymptotically equivalent to the expected information.

To prove that $J(\boldsymbol{\theta}) = \text{Var}[\nabla L(\boldsymbol{\theta})]$, suppose the data has density $f(\mathbf{y} \mid \boldsymbol{\theta})$ relative to some measure ν , which is usually ordinary volume measure or a discrete counting measure. We first note that the score conveniently has vanishing expectation because

$$\begin{aligned} \text{E}[\nabla L(\boldsymbol{\theta})] &= \int \frac{\nabla f(\mathbf{y} \mid \boldsymbol{\theta})}{f(\mathbf{y} \mid \boldsymbol{\theta})} f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= \nabla \int f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) \end{aligned}$$

and $\int f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) = 1$. Here the interchange of differentiation and expectation must be proved, but we will not stop to do so. See the references [176, 218]. The formal calculation

$$\begin{aligned} \text{E}[-d^2 L(\boldsymbol{\theta})] &= - \int \left[\frac{d^2 f(\mathbf{y} \mid \boldsymbol{\theta})}{f(\mathbf{y} \mid \boldsymbol{\theta})} - \frac{\nabla f(\mathbf{y} \mid \boldsymbol{\theta}) d f(\mathbf{y} \mid \boldsymbol{\theta})}{f(\mathbf{y} \mid \boldsymbol{\theta})^2} \right] f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= -d^2 \int f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &\quad + \int \nabla L(\boldsymbol{\theta}) dL(\boldsymbol{\theta}) f(\mathbf{y} \mid \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= -\mathbf{0} + \text{E}[\nabla L(\boldsymbol{\theta}) dL(\boldsymbol{\theta})] \end{aligned}$$

then completes the verification.

The score and expected information simplify considerably for exponential families of densities [22, 43, 110, 146, 202]. Based on equation (9.1), the score and expected information can be expressed succinctly in terms of the mean vector $\boldsymbol{\mu}(\boldsymbol{\theta}) = \text{E}[h(\mathbf{y})]$ and the variance matrix $\boldsymbol{\Sigma}(\boldsymbol{\theta}) = \text{Var}[h(\mathbf{y})]$ of the sufficient statistic $h(\mathbf{y})$. Our point of departure in deriving these quantities is the identity

$$dL(\boldsymbol{\theta}) = d\beta(\boldsymbol{\theta}) + h(\mathbf{y})^* d\gamma(\boldsymbol{\theta}). \tag{10.8}$$

If $\gamma(\boldsymbol{\theta})$ is linear in $\boldsymbol{\theta}$, then $J(\boldsymbol{\theta}) = -d^2 L(\boldsymbol{\theta}) = -d^2 \beta(\boldsymbol{\theta})$, and scoring coincides with Newton's method. If in addition $J(\boldsymbol{\theta})$ is positive definite, then $L(\boldsymbol{\theta})$ is strictly concave and possesses at most a single local maximum, which is necessarily the global maximum.

For an exponential family, the fact that $\text{E}[\nabla L(\boldsymbol{\theta})] = \mathbf{0}$ can be restated as

$$d\beta(\boldsymbol{\theta}) + \boldsymbol{\mu}(\boldsymbol{\theta})^* d\gamma(\boldsymbol{\theta}) = \mathbf{0}^*. \tag{10.9}$$

Subtracting equation (10.9) from equation (10.8) yields the alternative representation

$$dL(\boldsymbol{\theta}) = [h(\mathbf{y}) - \boldsymbol{\mu}(\boldsymbol{\theta})]^* d\gamma(\boldsymbol{\theta}) \tag{10.10}$$

of the first differential. This representation implies that the expected information is

$$J(\boldsymbol{\theta}) = \text{Var}[\nabla L(\boldsymbol{\theta})] = d\boldsymbol{\gamma}(\boldsymbol{\theta})^* \Sigma(\boldsymbol{\theta}) d\boldsymbol{\gamma}(\boldsymbol{\theta}). \quad (10.11)$$

To eliminate $d\boldsymbol{\gamma}(\boldsymbol{\theta})$ in equations (10.10) and (10.11), note that

$$\begin{aligned} d\boldsymbol{\mu}(\boldsymbol{\theta}) &= \int h(\mathbf{y}) df(\mathbf{y} | \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= \int h(\mathbf{y}) dL(\boldsymbol{\theta}) f(\mathbf{y} | \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= \int h(\mathbf{y}) [h(\mathbf{y}) - \boldsymbol{\mu}(\boldsymbol{\theta})]^* d\boldsymbol{\gamma}(\boldsymbol{\theta}) f(\mathbf{y} | \boldsymbol{\theta}) d\nu(\mathbf{y}) \\ &= \int [h(\mathbf{y}) - \boldsymbol{\mu}(\boldsymbol{\theta})][h(\mathbf{y}) - \boldsymbol{\mu}(\boldsymbol{\theta})]^* f(\mathbf{y} | \boldsymbol{\theta}) d\nu(\mathbf{y}) d\boldsymbol{\gamma}(\boldsymbol{\theta}) \\ &= \Sigma(\boldsymbol{\theta}) d\boldsymbol{\gamma}(\boldsymbol{\theta}). \end{aligned}$$

When $\Sigma(\boldsymbol{\theta})$ is invertible, this calculation implies $d\boldsymbol{\gamma}(\boldsymbol{\theta}) = \Sigma(\boldsymbol{\theta})^{-1} d\boldsymbol{\mu}(\boldsymbol{\theta})$, which in view of equations (10.10) and (10.11) yields

$$dL(\boldsymbol{\theta}) = [h(\mathbf{y}) - \boldsymbol{\mu}(\boldsymbol{\theta})]^* \Sigma(\boldsymbol{\theta})^{-1} d\boldsymbol{\mu}(\boldsymbol{\theta}) \quad (10.12)$$

$$J(\boldsymbol{\theta}) = d\boldsymbol{\mu}(\boldsymbol{\theta})^* \Sigma(\boldsymbol{\theta})^{-1} d\boldsymbol{\mu}(\boldsymbol{\theta}). \quad (10.13)$$

One can verify these formulas directly for the normal, Poisson, exponential, and multinomial distributions studied in the previous section. In each instance the sufficient statistic for case i is just y_i .

Based on equations (9.1), (10.12), and (10.13), Table 10.2 displays the loglikelihood, score vector, and expected information matrix for some commonly applied exponential families. In this table, x represents a single observation from the binomial, Poisson, and exponential families. For the multinomial family with m categories, $\mathbf{x} = (x_1, \dots, x_m)^*$ gives the category counts. The quantity μ denotes the mean of x for the Poisson and exponential families. For the binomial family, we express the mean np as the product of the number of trials n and the success probability p per trial. A similar convention holds for the multinomial family.

The multinomial family deserves further comment. Straightforward calculation shows that the variance matrix $\Sigma(\boldsymbol{\theta})$ has entries

$$n[1_{\{i=j\}} p_i(\boldsymbol{\theta}) - p_i(\boldsymbol{\theta}) p_j(\boldsymbol{\theta})].$$

Here the matrix $\Sigma(\boldsymbol{\theta})$ is singular, so the generalized inverse applies in formulas (10.12) and (10.13). In this case it is easier to derive the expected information by taking the expectation of the observed information given in Sect. 10.5.

In the ABO allele frequency estimation problem studied in Chaps. 8 and 9, scoring can be implemented by taking as basic parameters p_A and p_B

TABLE 10.2. Score and information for some exponential families

Family	$L(\theta)$	$\nabla L(\theta)$	$J(\theta)$
Binomial	$x \ln \frac{p}{1-p} + n \ln(1-p)$	$\frac{x-np}{p(1-p)} \nabla p$	$\frac{n}{p(1-p)} \nabla p dp$
Multinomial	$\sum_i x_i \ln p_i$	$\sum_i \frac{x_i}{p_i} \nabla p_i$	$\sum_i \frac{n}{p_i} \nabla p_i dp_i$
Poisson	$-\mu + x \ln \mu$	$-\nabla \mu + \frac{x}{\mu} \nabla \mu$	$\frac{1}{\mu} \nabla \mu d\mu$
Exponential	$-\ln \mu - \frac{x}{\mu}$	$-\frac{1}{\mu} \nabla \mu + \frac{x}{\mu^2} \nabla \mu$	$\frac{1}{\mu^2} \nabla \mu d\mu$

and expressing $p_O = 1 - p_A - p_B$. Scoring then leads to the same maximum likelihood point $(\hat{p}_A, \hat{p}_B, \hat{p}_O) = (.2136, .0501, .7363)$ as the EM algorithm. The quicker convergence of scoring here—four iterations as opposed to five starting from $(.3, .2, .5)$ —is often more dramatic in other problems. Scoring also has the advantage over EM of immediately providing asymptotic standard deviations of the parameter estimates. These are $(.0135, .0068, .0145)$ for the estimates $(\hat{p}_A, \hat{p}_B, \hat{p}_O)$.

10.7 The Gauss-Newton Algorithm

Armed with our better understanding of scoring, let us revisit nonlinear regression. Suppose that the n independent responses y_1, \dots, y_n are normally distributed with means $\mu_i(\theta)$ and variances σ^2/w_i , where the w_i are known constants. To estimate the mean parameter vector θ and the variance parameter σ^2 by scoring, we first write the loglikelihood up to a constant as the function

$$L(\phi) = -\frac{n}{2} \ln \sigma^2 - \frac{1}{2\sigma^2} \sum_{i=1}^n w_i [y_i - \mu_i(\theta)]^2 = -\frac{n}{2} \ln \sigma^2 - \frac{f(\theta)}{\sigma^2}$$

of the parameters $\phi = (\theta^*, \sigma^2)^*$.

Straightforward differentiation yields the score

$$\nabla L(\phi) = \begin{pmatrix} \frac{1}{\sigma^2} \sum_{i=1}^n w_i [y_i - \mu_i(\theta)] \nabla \mu_i(\theta) \\ -\frac{n}{2\sigma^2} + \frac{1}{2\sigma^4} \sum_{i=1}^n w_i [y_i - \mu_i(\theta)]^2 \end{pmatrix}.$$

To derive the expected information

$$J(\phi) = \begin{pmatrix} \frac{1}{\sigma^2} \sum_{i=1}^n w_i \nabla \mu_i(\theta) d\mu_i(\theta) & 0 \\ 0 & \frac{n}{2\sigma^4} \end{pmatrix},$$

we note that the observed information matrix can be written as the symmetric block matrix

$$-d^2L(\phi) = \begin{pmatrix} \mathbf{H}_{\theta\theta} & \mathbf{H}_{\theta\sigma^2} \\ \mathbf{H}_{\sigma^2\theta} & H_{\sigma^2\sigma^2} \end{pmatrix}.$$

The upper-left block $\mathbf{H}_{\theta\theta}$ equals $d^2f(\boldsymbol{\theta})/\sigma^2$ with $d^2f(\boldsymbol{\theta})$ given by equation (10.7). The displayed value of the expectation $E(\mathbf{H}_{\theta\theta})$ follows directly from the identity $E[y_i - \mu_i(\boldsymbol{\theta})] = 0$. The upper-right block $\mathbf{H}_{\theta\sigma^2}$ amounts to $-\nabla f(\boldsymbol{\theta})/\sigma^4$, and its expectation vanishes because again $E[y_i - \mu_i(\boldsymbol{\theta})] = 0$. Finally, the lower-right block $H_{\sigma^2\sigma^2}$ equals

$$-\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta})]^2.$$

Its expectation

$$E(H_{\sigma^2\sigma^2}) = -\frac{n}{2\sigma^4} + \frac{1}{\sigma^6} \sum_{i=1}^n w_i E\{[y_i - \mu_i(\boldsymbol{\theta})]^2\} = \frac{n}{2\sigma^4}$$

because $\text{Var}(y_i) = \sigma^2/w_i$. Readers experienced in calculating variances and covariances can verify the blocks of $J(\boldsymbol{\theta})$ by forming $\text{Var}[\nabla L(\boldsymbol{\theta})]$.

In any event, scoring updates $\boldsymbol{\theta}$ by

$$\begin{aligned} \boldsymbol{\theta}_{m+1} & \qquad \qquad \qquad (10.14) \\ &= \boldsymbol{\theta}_m + \left[\sum_{i=1}^n w_i \nabla \mu_i(\boldsymbol{\theta}_m) d\mu(\boldsymbol{\theta}_m) \right]^{-1} \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta}_m)] \nabla \mu_i(\boldsymbol{\theta}_m) \end{aligned}$$

and σ^2 by

$$\sigma_{m+1}^2 = \frac{1}{n} \sum_{i=1}^n w_i [y_i - \mu_i(\boldsymbol{\theta}_m)]^2.$$

The scoring algorithm (10.14) for $\boldsymbol{\theta}$ amounts to nothing more than the Gauss-Newton algorithm. The Gauss-Newton updates can be carried out blithely neglecting the updates of σ^2 .

10.8 Generalized Linear Models

The generalized linear model [202] deals with exponential families (9.1) in which the sufficient statistic $h(y)$ is y and the mean $\mu(\boldsymbol{\theta})$ of y completely determines the distribution of y . In many applications it is natural to postulate that $\mu(\boldsymbol{\theta}) = q(\mathbf{x}^* \boldsymbol{\theta})$ is a monotone function $q(s)$ of some linear combination of known predictors \mathbf{x} . The inverse of $q(s)$ is called the

TABLE 10.3. AIDS data from Australia during 1983–1986

Quarter	Deaths	Quarter	Deaths	Quarter	Deaths
1	0	6	4	11	20
2	1	7	9	12	25
3	2	8	18	13	37
4	3	9	23	14	45
5	1	10	31		

link function. In this setting, $d\mu(\boldsymbol{\theta}) = q'(\mathbf{x}^*\boldsymbol{\theta})\mathbf{x}^*$. It follows from equations (10.12) and (10.13) that if y_1, \dots, y_j are independent responses with corresponding predictor vectors $\mathbf{x}_1, \dots, \mathbf{x}_j$, then the score and expected information can be written as

$$\begin{aligned}\nabla L(\boldsymbol{\theta}) &= \sum_{i=1}^j \frac{y_i - \mu_i(\boldsymbol{\theta})}{\sigma_i^2(\boldsymbol{\theta})} q'(\mathbf{x}_i^*\boldsymbol{\theta})\mathbf{x}_i \\ J(\boldsymbol{\theta}) &= \sum_{i=1}^j \frac{1}{\sigma_i^2(\boldsymbol{\theta})} q'(\mathbf{x}_i^*\boldsymbol{\theta})^2 \mathbf{x}_i \mathbf{x}_i^*,\end{aligned}$$

where $\sigma_i^2(\boldsymbol{\theta}) = \text{Var}(y_i)$.

Table 10.3 contains quarterly data on AIDS deaths in Australia that illustrate the application of a generalized linear model [73, 273]. A simple plot of the data suggests exponential growth. A plausible model therefore involves Poisson distributed observations y_i with means $\mu_i(\boldsymbol{\theta}) = e^{\theta_1 + i\theta_2}$. Because this parameterization renders scoring equivalent to Newton's method, scoring gives the quick convergence noted in Table 10.4.

10.9 Accelerated MM

We now consider the question of how to accelerate the often excruciatingly slow convergence of the MM algorithm. The simplest device is to just double each MM step [61, 163]. Thus, if $F(\mathbf{x}_m)$ is the MM algorithm map from \mathbb{R}^p to \mathbb{R}^p , then we move to $\mathbf{x}_m + 2[F(\mathbf{x}_m) - \mathbf{x}_m]$ rather than to $F(\mathbf{x}_m)$. Step doubling is a standard tactic that usually halves the number of iterations until convergence. However, in many problems something more radical is necessary. Because Newton's method enjoys exceptionally quick convergence in a neighborhood of the optimal point, an attractive strategy is to amend the MM algorithm so that it resembles Newton's method. The papers [144, 145, 164] take up this theme from the perspective of optimizing the objective function by Newton's method. It is also possible to apply Newton's method to find a root of the equation $\mathbf{0} = \mathbf{x} - F(\mathbf{x})$. This alternative perspective has the advantage of dealing directly with the iterates of

TABLE 10.4. Scoring iterates for the AIDS model

Iteration	Step halves	θ_1	θ_2
1	0	0.0000	0.0000
2	3	-1.3077	0.4184
3	0	0.6456	0.2401
4	0	0.3744	0.2542
5	0	0.3400	0.2565
6	0	0.3396	0.2565

the MM algorithm. Let $G(\mathbf{x})$ denote the difference $G(\mathbf{x}) = \mathbf{x} - F(\mathbf{x})$. Because $G(\mathbf{x})$ has differential $dG(\mathbf{x}) = \mathbf{I} - dF(\mathbf{x})$, Newton's method iterates according to

$$\begin{aligned} \mathbf{x}_{m+1} &= \mathbf{x}_m - dG(\mathbf{x}_m)^{-1}G(\mathbf{x}_m) \\ &= \mathbf{x}_m - [\mathbf{I} - dF(\mathbf{x}_m)]^{-1}G(\mathbf{x}_m). \end{aligned} \quad (10.15)$$

If we can approximate $dF(\mathbf{x}_m)$ by a low-rank matrix \mathbf{M} , then we can replace $\mathbf{I} - dF(\mathbf{x}_m)$ by $\mathbf{I} - \mathbf{M}$ and explicitly form the inverse $(\mathbf{I} - \mathbf{M})^{-1}$. Let us see where this strategy leads.

Quasi-Newton methods operate by secant approximations [27, 32]. It is easy to generate a secant condition by taking two MM iterates starting from the current point \mathbf{x}_m . Close to the optimal point \mathbf{y} , the linear approximation

$$F \circ F(\mathbf{x}_m) - F(\mathbf{x}_m) \approx \mathbf{M}[F(\mathbf{x}_m) - \mathbf{x}_m]$$

holds, where $\mathbf{M} = dF(\mathbf{y})$. If \mathbf{v} is the vector $F \circ F(\mathbf{x}_m) - F(\mathbf{x}_m)$ and \mathbf{u} is the vector $F(\mathbf{x}_m) - \mathbf{x}_m$, then the secant condition is $\mathbf{M}\mathbf{u} = \mathbf{v}$. In fact, the best results may require several secant conditions $\mathbf{M}\mathbf{u}_i = \mathbf{v}_i$ for $i = 1, \dots, q$, where $q \leq p$. These can be generated at the current iterate \mathbf{x}_m and the previous $q - 1$ iterates. For convenience represent the secant conditions in the matrix form $\mathbf{M}\mathbf{U} = \mathbf{V}$ for $\mathbf{U} = (\mathbf{u}_1, \dots, \mathbf{u}_q)$ and $\mathbf{V} = (\mathbf{v}_1, \dots, \mathbf{v}_q)$. Example 5.2.7 shows that the choice $\mathbf{M} = \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*$ minimizes the Frobenius norm of \mathbf{M} subject to the secant constraint $\mathbf{M}\mathbf{U} = \mathbf{V}$. In practice, it is better to make a controlled approximation to $dF(\mathbf{y})$ than a wild guess.

To apply the approximation, we must invert the matrix $\mathbf{I} - \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*$. Fortunately, we have the explicit inverse

$$[\mathbf{I} - \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*]^{-1} = \mathbf{I} + \mathbf{V}[\mathbf{U}^*\mathbf{U} - \mathbf{U}^*\mathbf{V}]^{-1}\mathbf{U}^*. \quad (10.16)$$

The reader can readily check this variant of the Sherman-Morrison formula. It is noteworthy that the $q \times q$ matrix $\mathbf{U}^*\mathbf{U} - \mathbf{U}^*\mathbf{V}$ is trivial to invert for q small even when p is large. With these results in hand, the Newton update

(10.15) can be replaced by the quasi-Newton update

$$\begin{aligned}\mathbf{x}_{m+1} &= \mathbf{x}_m - [\mathbf{I} - \mathbf{V}(\mathbf{U}^*\mathbf{U})^{-1}\mathbf{U}^*]^{-1}[\mathbf{x}_m - F(\mathbf{x}_m)] \\ &= \mathbf{x}_m - [\mathbf{I} + \mathbf{V}(\mathbf{U}^*\mathbf{U} - \mathbf{U}^*\mathbf{V})^{-1}\mathbf{U}^*][\mathbf{x}_m - F(\mathbf{x}_m)] \\ &= F(\mathbf{x}_m) - \mathbf{V}(\mathbf{U}^*\mathbf{U} - \mathbf{U}^*\mathbf{V})^{-1}\mathbf{U}^*[\mathbf{x}_m - F(\mathbf{x}_m)].\end{aligned}$$

The special case $q = 1$ is interesting in its own right. A brief calculation shows that the quasi-Newton update for $q = 1$ is

$$\begin{aligned}\mathbf{x}_{m+1} &= (1 - c_m)F(\mathbf{x}_m) + c_m F \circ F(\mathbf{x}_m) \\ c_m &= \frac{\|F(\mathbf{x}_m) - \mathbf{x}_m\|^2}{[F \circ F(\mathbf{x}_m) - 2F(\mathbf{x}_m) + \mathbf{x}_m]^* [F(\mathbf{x}_m) - \mathbf{x}_m]}.\end{aligned}\tag{10.17}$$

This quasi-Newton acceleration enjoys several desirable properties in high-dimensional problems. First, the computational effort per iteration is relatively light: two MM updates and a few matrix times vector multiplications. Second, memory demands are also light. If we fix q in advance, the most onerous requirement is storage of the secant matrices \mathbf{U} and \mathbf{V} . These two matrices can be updated by replacing the earliest retained secant pair by the latest secant pair generated. Third, the whole scheme is consistent with linear constraints. Thus, if the parameter space satisfies a linear constraint $\mathbf{w}^*\mathbf{x} = a$ for all feasible \mathbf{x} , then the quasi-Newton iterates also satisfy $\mathbf{w}^*\mathbf{x}_m = a$ for all m . This claim follows from the equalities $\mathbf{w}^*F(\mathbf{x}) = a$ and $\mathbf{w}^*\mathbf{V} = \mathbf{0}$. Finally, if the quasi-Newton update at \mathbf{x}_m fails the ascent or descent test, then one can always revert to the second MM update $F \circ F(\mathbf{x}_m)$. Balanced against these advantages is the failure of the quasi-Newton acceleration to respect parameter lower and upper bounds.

Example 10.9.1 *A Mixture of Poissons*

Problem 9 of Chap. 9 describes a Poisson mixture model for mortality data from *The London Times*. Starting from the method of moments estimates $(\mu_{01}, \mu_{02}, \pi_0) = (1.101, 2.582, .2870)$, the EM algorithm takes an excruciating 535 iterations for the loglikelihood $L(\theta)$ to attain its maximum of -1989.946 . Even worse, it takes 1,749 iterations for the parameters to reach the maximum likelihood estimates $(\hat{\mu}_1, \hat{\mu}_2, \hat{\pi}) = (1.256, 2.663, .3599)$. The sizable difference in convergence rates to the maximum loglikelihood and the maximum likelihood estimates indicates that the likelihood surface is quite flat. In contrast, the accelerated EM algorithm converges to the maximum loglikelihood in about 10–150 iterations, depending on the value of q . Figure 10.1 plots the progress of the EM algorithm and the different versions of the quasi-Newton acceleration. Titterington et al. [259] report that Newton's method typically takes 8–11 iterations to converge when it converges for these data. For about a third of their initial points, Newton's method fails. ■

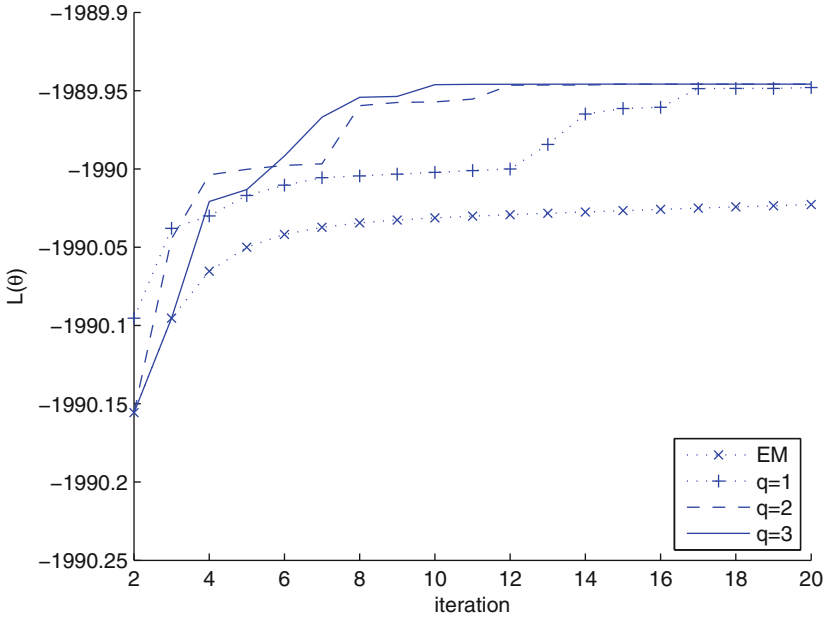


FIGURE 10.1. MM acceleration for the mixture of Poissons example

Although we have couched quasi-Newton acceleration in terms of the MM algorithm, it applies to any optimization algorithm with reasonably smooth parameter updates. Block relaxation is another potential beneficiary. Block relaxation shares the ascent-descent property of the MM algorithm, so if acceleration fails to improve the objective function, then one can still make progress by reverting to the original double step executed in constructing a new secant.

10.10 Problems

1. What happens when you apply Newton's method to the functions

$$f(x) = \begin{cases} \sqrt{x} & x \geq 0 \\ -\sqrt{-x} & x < 0 \end{cases}$$

and $g(x) = \sqrt[3]{x}$?

2. Consider a function $f(x) = (x - r)^k g(x)$ with a root r of multiplicity k . If $g'(x)$ is continuous at r , and the Newton iterates x_m converge to r , then show that the iterates satisfy

$$\lim_{m \rightarrow \infty} \frac{|x_{m+1} - r|}{|x_m - r|} = 1 - \frac{1}{k}.$$

3. As an illustration of Problem 2, use Newton's method to extract a root of the polynomials $p_1(x) = x^2 - 1$ and $p_2(x) = x^2 - 2x + 1$ starting from $x_0 = 2$. Notice how much more slowly convergence occurs for $p_2(x)$ than for $p_1(x)$.
4. Suppose the real-valued $f(x)$ satisfies $f'(x) > 0$ and $f''(x) > 0$ for all x in its domain (d, ∞) . If the equation $f(x) = 0$ has a root r , then demonstrate that r is unique and that Newton's method converges to r regardless of its starting point. Further, prove that x_m converges monotonically to r from above when $x_0 > r$ and that $x_1 > r$ when $x_0 < r$. How are these results pertinent to Example 10.2.2?
5. Problem 4 applies to polynomials $p(x)$ having only real roots. Suppose $p(x)$ is a polynomial of degree d with roots $r_1 < r_2 < \dots < r_d$ and leading coefficient $c_d > 0$. Show that on the interval (r_d, ∞) the functions $p(x)$, $p'(x)$, and $p''(x)$ are all positive. Hence, if we seek r_d by Newton's method starting at $x_0 > r_d$, then the iterates x_m decrease monotonically to r_d . (Hint: According to Rolle's theorem, what can we say about the roots of $p'(x)$ and $p''(x)$?)
6. Suppose that the polynomial $p(x)$ has the known roots r_1, \dots, r_d . Maehly's algorithm [246] attempts to extract one additional root r_{d+1} by iterating via

$$x_{m+1} = x_m - \frac{p(x_m)}{p'(x_m) - \sum_{k=1}^d \frac{p(x_m)}{x_m - r_k}}.$$

Show that this is just a disguised version of Newton's method. It has the virtue of being more numerically accurate than Newton's method applied to the deflated polynomial calculated from $p(x)$ by synthetic division. (Hint: Consider the polynomial $q(x) = p(x) \prod_{k=1}^d (x - r_k)^{-1}$.)

7. Apply Maehly's algorithm as sketched in Problem 6 to find the roots of the polynomial $p(x) = x^4 - 12x^3 + 47x^2 - 60x$.
8. Consider the map

$$f(\mathbf{x}) = \begin{pmatrix} x_1^2 + x_2^2 - 2 \\ x_1 - x_2 \end{pmatrix}$$

of the plane into itself. Show that $f(\mathbf{x}) = \mathbf{0}$ has the roots $-\mathbf{1}$ and $\mathbf{1}$ and no other roots. Prove that Newton's method iterates according to

$$x_{m+1,1} = x_{m+1,2} = \frac{x_{m1}^2 + x_{m2}^2 + 2}{2(x_{m1} + x_{m2})}$$

and that these iterates converge to the root $-\mathbf{1}$ if $x_{01} + x_{02}$ is negative and to the root $\mathbf{1}$ if $x_{01} + x_{02}$ is positive. If $x_{01} + x_{02} = 0$, then the

first iterate is undefined. Finally, prove that

$$\lim_{m \rightarrow \infty} \frac{|x_{m+1,1} - y_1|}{|x_{m1} - y_1|^2} = \lim_{m \rightarrow \infty} \frac{|x_{m+1,2} - y_2|}{|x_{m2} - y_2|^2} = \frac{1}{2},$$

where \mathbf{y} is the root relevant to the initial point \mathbf{x}_0 .

9. Continuing Example 10.2.5, consider iterating according to

$$\mathbf{B}_{m+1} = \mathbf{B}_m \sum_{i=0}^j (\mathbf{I} - \mathbf{A}\mathbf{B}_m)^i \quad (10.18)$$

to find \mathbf{A}^{-1} [138]. Example 10.2.5 covers the special case $j = 1$. Verify the alternative representation

$$\mathbf{B}_{m+1} = \sum_{i=0}^j (\mathbf{I} - \mathbf{B}_m \mathbf{A})^i \mathbf{B}_m,$$

and use it to prove that \mathbf{B}_{m+1} is symmetric whenever \mathbf{A} and \mathbf{B}_m are. Also show that

$$\mathbf{A}^{-1} - \mathbf{B}_{m+1} = (\mathbf{A}^{-1} - \mathbf{B}_m)[\mathbf{A}(\mathbf{A}^{-1} - \mathbf{B}_m)]^j.$$

From this last identity deduce the norm inequality

$$\|\mathbf{A}^{-1} - \mathbf{B}_{m+1}\| \leq \|\mathbf{A}\|^j \|\mathbf{A}^{-1} - \mathbf{B}_m\|^{j+1}.$$

Thus, the algorithm converges at a cubic rate when $j = 2$, at a quartic rate when $j = 3$, and so forth.

10. Example 10.2.2 can be adapted to extract the n th root of a positive semidefinite matrix \mathbf{A} [107]. Consider the iteration scheme

$$\mathbf{B}_{m+1} = \frac{n-1}{n} \mathbf{B}_m + \frac{1}{n} \mathbf{B}_m^{-n+1} \mathbf{A}$$

starting with $\mathbf{B}_0 = c\mathbf{I}$ for some positive constant c . Show by induction that (a) \mathbf{B}_m commutes with \mathbf{A} , (b) \mathbf{B}_m is symmetric, and (c) \mathbf{B}_m is positive definite. To prove that \mathbf{B}_m converges to $\mathbf{A}^{1/n}$, consider the spectral decomposition $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^*$ of \mathbf{A} with \mathbf{D} diagonal and \mathbf{U} orthogonal. Show that \mathbf{B}_m has a similar spectral decomposition $\mathbf{B}_m = \mathbf{U}\mathbf{D}_m\mathbf{U}^*$ and that the i th diagonal entries of \mathbf{D}_m and \mathbf{D} satisfy

$$d_{m+1,i} = \frac{n-1}{n} d_{mi} + \frac{1}{n} d_{mi}^{-n+1} d_i.$$

Example 10.2.2 implies that d_{mi} converges to $\sqrt[n]{d_i}$ when $d_i > 0$. This convergence occurs at a fast quadratic rate as explained in Proposition 12.2.2. If $d_i = 0$, then d_{mi} converges to 0 at the linear rate $\frac{n-1}{n}$.

11. Program the algorithm of Problem 10 and extract the square roots of the two matrices

$$\begin{pmatrix} 1 & 1 \\ 1 & 1 \end{pmatrix}, \quad \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}.$$

Describe the apparent rate of convergence in each case and any difficulties you encounter with roundoff error.

12. Prove that the increment (10.5) can be expressed as

$$\begin{aligned} \Delta \mathbf{x}_m &= -\mathbf{H}_m^{-1/2} \left[\mathbf{I} - \mathbf{H}_m^{-1/2} \mathbf{V}^* (\mathbf{V} \mathbf{H}_m^{-1} \mathbf{V}^*)^{-1} \mathbf{V} \mathbf{H}_m^{-1/2} \right] \mathbf{H}_m^{-1/2} \nabla f(\mathbf{x}_m) \\ &= -\mathbf{H}_m^{-1/2} (\mathbf{I} - \mathbf{P}_m) \mathbf{H}_m^{-1/2} \nabla f(\mathbf{x}_m) \end{aligned}$$

using the symmetric square root $\mathbf{H}_m^{-1/2}$ of \mathbf{H}_m^{-1} . Check that the matrix \mathbf{P}_m is a projection in the sense that $\mathbf{P}_m^* = \mathbf{P}_m$ and $\mathbf{P}_m^2 = \mathbf{P}_m$ and that these properties carry over to $\mathbf{I} - \mathbf{P}_m$. Now argue that

$$-df(\mathbf{x}_m) \Delta \mathbf{x}_m = \|(\mathbf{I} - \mathbf{P}_m) \mathbf{H}_m^{-1/2} \nabla f(\mathbf{x}_m)\|^2$$

and consequently that step halving is bound to produce a decrease in $f(\mathbf{x})$ if $(\mathbf{I} - \mathbf{P}_m) \mathbf{H}_m^{-1/2} \nabla f(\mathbf{x}_m) \neq \mathbf{0}$.

13. Show that Newton's method converges in one iteration to the minimum of

$$f(\boldsymbol{\theta}) = \frac{1}{2} \boldsymbol{\theta}^* \mathbf{A} \boldsymbol{\theta} + \mathbf{b}^* \boldsymbol{\theta} + c$$

when the symmetric matrix \mathbf{A} is positive definite. Note that this implies that the Gauss-Newton algorithm (10.14) converges in a single step when the regression functions $\mu_i(\boldsymbol{\theta})$ are linear.

14. In Example 10.4.2, digamma and trigamma functions $\psi(t)$ and $\psi'(t)$ must be evaluated. Show that these functions satisfy the recurrence relations

$$\begin{aligned} \psi(t) &= -t^{-1} + \psi(t+1) \\ \psi'(t) &= t^{-2} + \psi'(t+1). \end{aligned}$$

Thus, if $\psi(t)$ and $\psi'(t)$ can be accurately evaluated via asymptotic expansions for large t , then they can be accurately evaluated for small t . For example, it is known that $\psi(t) = \ln t - (2t)^{-1} + O(t^{-2})$ and $\psi'(t) = t^{-1} + (\sqrt{2}t)^{-2} + O(t^{-3})$ as $t \rightarrow \infty$.

15. Compute the score vector and the observed and expected information matrices for the Dirichlet distribution (10.6). Explicitly invert the expected information using the Sherman-Morrison formula.

16. Verify the score and information entries in Table 10.2.
17. Let $g(x)$ and $h(x)$ be probability densities defined on the real line. Show that the admixture density $f(x) = \theta g(x) + (1 - \theta)h(x)$ for $\theta \in [0, 1]$ has score and expected information

$$L'(\theta) = \frac{g(x) - h(x)}{\theta g(x) + (1 - \theta)h(x)}$$

$$J(\theta) = \int \frac{[g(x) - h(x)]^2}{\theta g(x) + (1 - \theta)h(x)} dx$$

$$= \frac{1}{\theta(1 - \theta)} \left[1 - \int \frac{g(x)h(x)}{\theta g(x) + (1 - \theta)h(x)} dx \right].$$

What happens to $J(\theta)$ when $g(x)$ and $h(x)$ coincide? What does $J(\theta)$ equal when $g(x)$ and $h(x)$ have nonoverlapping domains? (Hint: The identities

$$h - g = \frac{\theta g + (1 - \theta)h - g}{1 - \theta}, \quad g - h = \frac{\theta g + (1 - \theta)h - h}{\theta}$$

will help.)

18. A quantal response model involves independent binomial observations y_1, \dots, y_j with n_i trials and success probability $\pi_i(\boldsymbol{\theta})$ per trial for the i th observation. If \mathbf{x}_i is a predictor vector and $\boldsymbol{\theta}$ a parameter vector, then the specification

$$\pi_i(\boldsymbol{\theta}) = \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}}}$$

gives a generalized linear model. Use the scoring algorithm to estimate $\boldsymbol{\theta} = (-5.1316, 0.0677)^*$ for the ingot data of Cox [53] displayed in Table 10.5.

TABLE 10.5. Ingot data for a quantal response model

Trials n_i	Observation y_i	Covariate x_{i1}	Covariate x_{i2}
55	0	1	7
157	2	1	14
159	7	1	27
16	3	1	57

19. In robust regression it is useful to consider location-scale families with densities of the form

$$\frac{c}{\sigma} e^{-\rho\left(\frac{x-\mu}{\sigma}\right)}, \quad x \in (-\infty, \infty). \tag{10.19}$$

Here $\rho(r)$ is a strictly convex even function, decreasing to the left of 0 and symmetrically increasing to the right of 0. Without loss of generality, one can take $\rho(0) = 0$. The normalizing constant c is determined by $c \int_{-\infty}^{\infty} e^{-\rho(r)} dr = 1$. Show that a random variable X with density (10.19) has mean μ and variance

$$\text{Var}(X) = c\sigma^2 \int_{-\infty}^{\infty} r^2 e^{-\rho(r)} dr.$$

If μ depends on a parameter vector θ , demonstrate that the score corresponding to a single observation $X = x$ amounts to

$$\nabla L(\phi) = \begin{pmatrix} \frac{1}{\sigma} \rho'(\frac{x-\mu}{\sigma}) \nabla \mu \\ -\frac{1}{\sigma} + \rho'(\frac{x-\mu}{\sigma}) \frac{x-\mu}{\sigma^2} \end{pmatrix}$$

for $\phi = (\theta^*, \sigma)^*$. Finally, prove that the expected information $J(\phi)$ is block diagonal with upper-left block

$$\frac{c}{\sigma^2} \int_{-\infty}^{\infty} \rho''(r) e^{-\rho(r)} dr \nabla \mu(\theta) d\mu(\theta)$$

and lower-right block

$$\frac{c}{\sigma^2} \int_{-\infty}^{\infty} \rho''(r) r^2 e^{-\rho(r)} dr + \frac{1}{\sigma^2}.$$

20. In the context of Problem 19, take $\rho(r) = \ln \cosh^2(\frac{r}{2})$. Show that this corresponds to the logistic distribution with density

$$f(x) = \frac{e^{-x}}{(1 + e^{-x})^2}.$$

Compute the integrals

$$\begin{aligned} \frac{\pi^2}{3} &= c \int_{-\infty}^{\infty} r^2 e^{-\rho(r)} dr \\ \frac{1}{3} &= c \int_{-\infty}^{\infty} \rho''(r) e^{-\rho(r)} dr \\ \frac{1}{3} + \frac{\pi^2}{9} &= c \int_{-\infty}^{\infty} \rho''(r) r^2 e^{-\rho(r)} dr + 1 \end{aligned}$$

determining the variance and expected information of the density (10.19) for this choice of $\rho(r)$.

21. Continuing Problems 19 and 20, compute the normalizing constant c and the three integrals determining the variance and expected information for Huber's function

$$\rho(r) = \begin{cases} \frac{r^2}{2} & |r| \leq k \\ k|r| - \frac{k^2}{2} & |r| > k. \end{cases}$$

22. A family of discrete density functions $p_n(\theta)$ defined on $\{0, 1, \dots\}$ and indexed by a parameter $\theta > 0$ is said to be a power series family if for all n

$$p_n(\theta) = \frac{c_n \theta^n}{g(\theta)}, \quad (10.20)$$

where $c_n \geq 0$, and where $g(\theta) = \sum_{n=0}^{\infty} c_n \theta^n$ is the appropriate normalizing constant. If y_1, \dots, y_j are independent observations from the discrete density (10.20), then show that the maximum likelihood estimate of θ is a root of the equation

$$\frac{1}{j} \sum_{i=1}^j y_i = \frac{\theta g'(\theta)}{g(\theta)}. \quad (10.21)$$

Prove that the expected information in a single observation is

$$J(\theta) = \frac{\sigma^2(\theta)}{\theta^2},$$

where $\sigma^2(\theta)$ is the variance of the density (10.20).

23. Continuing problem 22, equation (10.21) suggests that one can find the maximum likelihood estimate $\hat{\theta}$ by iterating via

$$\theta_{m+1} = \frac{\bar{x}g(\theta_m)}{g'(\theta_m)} = f(\theta_m),$$

where \bar{x} is the sample mean. The question now arises whether this iteration scheme is likely to converge to $\hat{\theta}$. Local convergence hinges on the condition $|f'(\hat{\theta})| < 1$. When this condition is true, the map $\theta_{m+1} = f(\theta_m)$ is locally contractive near the fixed point $\hat{\theta}$. Prove that

$$f'(\hat{\theta}) = 1 - \frac{\sigma^2(\hat{\theta})}{\mu(\hat{\theta})},$$

where

$$\mu(\theta) = \frac{\theta g'(\theta)}{g(\theta)}$$

is the mean of a single realization. Thus, convergence depends on the ratio of the variance to the mean. (Hints: By differentiating $g(\theta)$ it is easy to compute the mean and the second factorial moment

$$E[X(X-1)] = \frac{\theta^2 g''(\theta)}{g(\theta)}.$$

Substitute this in $f'(\hat{\theta})$, recall $\text{Var}(X) = E[X(X-1)] + E(X) - E(X)^2$, and invoke equality (10.21).)

24. In the Gauss-Newton algorithm (10.14), the matrix

$$d^2 f(\boldsymbol{\theta}_m) = \sum_{i=1}^j w_i \nabla \mu_i(\boldsymbol{\theta}_m) d\mu(\boldsymbol{\theta}_m)$$

can be singular or nearly so. To cure this ill, Marquardt suggested choosing $\lambda > 0$, substituting

$$\mathbf{H}_m = \sum_{i=1}^j w_i \nabla \mu_i(\boldsymbol{\theta}_m) d\mu(\boldsymbol{\theta}_m) + \lambda I$$

for $d^2 f(\boldsymbol{\theta}_m)$, and iterating according to

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + \mathbf{H}_m^{-1} \sum_{i=1}^j w_i [x_i - \mu_i(\boldsymbol{\theta}_m)] \nabla \mu_i(\boldsymbol{\theta}_m). \quad (10.22)$$

Prove that the increment $\Delta \boldsymbol{\theta}_m = \boldsymbol{\theta}_{m+1} - \boldsymbol{\theta}_m$ proposed in equation (10.22) minimizes the criterion

$$\frac{1}{2} \sum_{i=1}^j w_i [x_i - \mu_i(\boldsymbol{\theta}_m) - d\mu_i(\boldsymbol{\theta}_m) \Delta \boldsymbol{\theta}_m]^2 + \frac{\lambda}{2} \|\Delta \boldsymbol{\theta}_m\|^2.$$

25. Survival analysis deals with nonnegative random variables T modeling random lifetimes. Let such a random variable $T \geq 0$ have density function $f(t)$ and distribution function $F(t)$. The hazard function

$$\begin{aligned} h(t) &= \lim_{s \downarrow 0} \frac{\Pr(t < T \leq t + s \mid T > t)}{s} \\ &= \frac{f(t)}{1 - F(t)} \end{aligned}$$

represents the instantaneous rate of death under lifetime T . Statisticians call the right-tail probability $1 - F(t) = S(t)$ the survival function and view $h(t)$ as the derivative

$$h(t) = -\frac{d}{dt} \ln S(t).$$

The cumulative hazard function $H(t) = \int_0^t h(s) ds$ obviously satisfies the identity

$$S(t) = e^{-H(t)}.$$

In Cox's proportional hazards model, longevity depends not only on time but also predictors. This is formalized by taking

$$h(t) = \lambda(t) e^{\mathbf{x}^* \boldsymbol{\alpha}},$$

where \mathbf{x} and $\boldsymbol{\alpha}$ are column vectors of predictors and regression coefficients, respectively. For instance, \mathbf{x} might be $(1, d)^*$, where d indicates dosage of a life-prolonging drug.

Many clinical trials involve right censoring. In other words, instead of observing a lifetime $T = t$, we observe $T > t$. Censored and ordinary data can be mixed in the same study. Generally, each observation T comes with a censoring indicator W . If T is censored, then $W = 1$; otherwise, $W = 0$.

(a) Show that

$$H(t) = \Lambda(t)e^{\mathbf{x}^* \boldsymbol{\alpha}},$$

where

$$\Lambda(t) = \int_0^t \lambda(s) ds.$$

In the Weibull proportional hazards model, $\lambda(t) = \beta t^{\beta-1}$. Show that this translates into the survival and density functions

$$\begin{aligned} S(t) &= e^{-t^\beta e^{\mathbf{x}^* \boldsymbol{\alpha}}} \\ f(t) &= \beta t^{\beta-1} e^{\mathbf{x}^* \boldsymbol{\alpha}} e^{-t^\beta e^{\mathbf{x}^* \boldsymbol{\alpha}}}. \end{aligned}$$

(b) Consider n independent possibly censored observations t_1, \dots, t_n with corresponding predictor vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ and censoring indicators w_1, \dots, w_n . Prove that the loglikelihood of the data is

$$L(\boldsymbol{\alpha}, \beta) = \sum_{i=1}^n w_i \ln S_i(t_i) + \sum_{i=1}^n (1 - w_i) \ln f_i(t_i),$$

where $S_i(t)$ and $f_i(t)$ are the survival and density functions of the i th case.

(c) Calculate the score and observed information for the Weibull model as posed. The observed information is

$$\begin{aligned} -d^2 L(\boldsymbol{\alpha}, \beta) &= \sum_{i=1}^n t_i^\beta e^{\mathbf{x}_i^* \boldsymbol{\alpha}} \begin{pmatrix} \mathbf{x}_i \\ \ln t_i \end{pmatrix} \begin{pmatrix} \mathbf{x}_i \\ \ln t_i \end{pmatrix}^* \\ &\quad + \sum_{i=1}^n (1 - w_i) \begin{pmatrix} 0 & 0 \\ 0 & \beta^{-2} \end{pmatrix}. \end{aligned}$$

(d) Show that the loglikelihood $L(\boldsymbol{\alpha}, \beta)$ for the Weibull model is concave. Demonstrate that it is strictly concave if and only if the n vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ span \mathbb{R}^m , where $\boldsymbol{\alpha}$ has m components.

(e) Describe in detail how you would implement Newton's method for finding the maximum likelihood estimate of the parameter vector (α, β) . What difficulties might you encounter? Why is concavity of the loglikelihood helpful?

26. Write a computer program and reproduce the iterates displayed in Table 10.4.

27. Let x_1, \dots, x_m be a random sample from the gamma density

$$f(x) = \Gamma(\alpha)^{-1} \beta^\alpha x^{\alpha-1} e^{-\beta x}$$

on $(0, \infty)$. Find the score, observed information, and expected information for the parameters α and β , and demonstrate that Newton's method and scoring coincide.

28. Continuing Problem 27, derive the method of moments estimators

$$\hat{\alpha} = \frac{\bar{x}^2}{s^2}, \quad \hat{\beta} = \frac{\bar{x}}{s^2},$$

where $\bar{x} = \frac{1}{m} \sum_{i=1}^m x_i$ and $s^2 = \frac{1}{m} \sum_{i=1}^m (x_i - \bar{x})^2$ are the sample mean and variance, respectively. These are not necessarily the best explicit estimators of the two parameters. Show that setting the score function equal to $\mathbf{0}$ implies that $\beta = \alpha/\bar{x}$ is a stationary point of the loglikelihood $L(\alpha, \beta)$ of the sample x_1, \dots, x_m for α fixed. Why does $\beta = \alpha/\bar{x}$ furnish the maximum? Now argue that substituting this value of β in the loglikelihood reduces maximum likelihood estimation to optimization of the profile loglikelihood

$$L(\alpha) = m\alpha \ln \alpha - m\alpha \ln \bar{x} - m \ln \Gamma(\alpha) + m(\alpha - 1) \overline{\ln x} - m\alpha.$$

Here $\overline{\ln x} = \frac{1}{m} \sum_{i=1}^m \ln x_i$. There are two nasty terms in $L(\alpha)$. One is $\alpha \ln \alpha$, and the other is $\ln \Gamma(\alpha)$. We can eliminate both by appealing to a version of Stirling's formula. Ordinarily Stirling's formula is only applied for large factorials. This limitation is inconsistent with small α . However, Gosper's version of Stirling's formula is accurate for all arguments. This little-known version of Stirling's formula says that

$$\Gamma(\alpha + 1) \approx \sqrt{(\alpha + 1/6)2\pi} \alpha^\alpha e^{-\alpha}.$$

Given that $\Gamma(\alpha) = \Gamma(\alpha + 1)/\alpha$, show that the application of Gosper's formula leads to the approximate maximum likelihood estimate

$$\hat{\alpha} = \frac{3 - d + \sqrt{(3 - d)^2 + 24d}}{12d},$$

where $d = \ln \bar{x} - \overline{\ln x}$ [47]. Why is this estimate of α positive? Why does one take the larger root of the defining quadratic?

29. In the multilogit model, items are drawn from m categories. Let y_i denote the i th outcome of l independent draws and \mathbf{x}_i a corresponding predictor vector. The probability π_{ij} that $y_i = j$ is given by

$$\pi_{ij}(\boldsymbol{\theta}) = \begin{cases} \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}_j}}{1 + \sum_{k=1}^{m-1} e^{\mathbf{x}_i^* \boldsymbol{\theta}_k}} & 1 \leq j < m \\ \frac{1}{1 + \sum_{k=1}^{m-1} e^{\mathbf{x}_i^* \boldsymbol{\theta}_k}} & j = m. \end{cases}$$

Find the loglikelihood, score, observed information, and expected information. Demonstrate that Newton's method and scoring coincide. (Hint: You can achieve compact expressions by stacking vectors and using matrix Kronecker products.)

30. Derive formulas (10.16) and (10.17).

11

Conjugate Gradient and Quasi-Newton

11.1 Introduction

Our discussion of Newton's method has highlighted both its strengths and its weaknesses. Related algorithms such as scoring and Gauss-Newton exploit special features of the objective function $f(\mathbf{x})$ in overcoming the defects of Newton's method. We now consider algorithms that apply to generic functions $f(\mathbf{x})$. These algorithms also operate by locally approximating $f(\mathbf{x})$ by a strictly convex quadratic function. Indeed, the guiding philosophy behind many modern optimization algorithms is to see what techniques work well with quadratic functions and then to modify the best techniques to accommodate generic functions.

The conjugate gradient algorithm [94, 127] is noteworthy for three properties: (a) it minimizes a quadratic function $f(\mathbf{x})$ from \mathbb{R}^n to \mathbb{R} in n steps, (b) it does not require evaluation of $d^2f(\mathbf{x})$, and (c) it does not involve storage or inversion of any $n \times n$ matrices. Property (c) makes the method particularly suitable for optimization in high-dimensional settings. One of the drawbacks of the conjugate gradient method is that it requires exact line searches.

Quasi-Newton algorithms [10, 56, 91, 93] enjoy properties (a) and (b) but not property (c). In compensation for the failure of (c), inexact line searches are usually adequate with quasi-Newton algorithms. Furthermore, quasi-Newton methods adapt more readily to parameter constraints. Except for a discussion of trust regions, the current chapter considers only unconstrained optimization problems.

11.2 Centers of Spheres and Centers of Ellipsoids

As an introduction to many of the central ideas of the chapter, it is instructive to explore a simple algorithm for finding the center of a sphere. The fact that we already know the answer should not deter us from considering algorithmic issues. If the center is the origin, then obviously we can find it by minimizing the scaled distance function

$$g(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2 = \frac{1}{2}\sum_{i=1}^n x_i^2$$

with gradient $\nabla g(\mathbf{x}) = \mathbf{x}$. In cyclic coordinate descent, we minimize $g(\mathbf{x})$ along each coordinate direction in turn, starting from a point \mathbf{x}_1 and generating successive points $\mathbf{x}_2, \dots, \mathbf{x}_{n+1}$. The search along coordinate direction \mathbf{e}_i at iteration i amounts to minimizing the function

$$g(\mathbf{x}_i + t\mathbf{e}_i) = \frac{1}{2}(x_{ii} + t)^2 + \frac{1}{2}\sum_{j \neq i} x_{ij}^2$$

of the scalar t . The minimum occurs at $t = -x_{ii}$ and yields \mathbf{x}_{i+1} . It is trivial to check that this procedure achieves the minimum in n iterations and satisfies at iteration i the identities

$$\mathbf{e}_i^* \mathbf{e}_j = 0 \quad \text{and} \quad dg(\mathbf{x}_i)\mathbf{e}_j = 0 \quad (11.1)$$

for all $j < i$. Furthermore, \mathbf{x}_{i+1} minimizes the function

$$h(t_1, \dots, t_i) = g\left(\mathbf{x}_1 + \sum_{j=1}^i t_j \mathbf{e}_j\right)$$

defined on the i -dimensional plane $\mathbf{x}_1 + t_1\mathbf{e}_1 + \dots + t_i\mathbf{e}_i$ formed from all linear combinations of the first i search directions. Because of the spherical symmetry of the function $g(\mathbf{x})$, we can substitute any set of nontrivial orthogonal vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ and reach the same conclusions.

If we consider an arbitrary strictly convex quadratic function

$$\begin{aligned} f(\mathbf{y}) &= \frac{1}{2}\mathbf{y}^* \mathbf{A} \mathbf{y} + \mathbf{b}^* \mathbf{y} + c \\ &= \frac{1}{2}(\mathbf{y} + \mathbf{A}^{-1}\mathbf{b})^* \mathbf{A} (\mathbf{y} + \mathbf{A}^{-1}\mathbf{b}) - \frac{1}{2}\mathbf{b}^* \mathbf{A}^{-1}\mathbf{b} + c, \end{aligned} \quad (11.2)$$

then the situation becomes more interesting. Here the matrix \mathbf{A} is positive definite, so there is no doubt that its inverse \mathbf{A}^{-1} exists. Because the minimum of $f(\mathbf{y})$ occurs at $\mathbf{y} = -\mathbf{A}^{-1}\mathbf{b}$, any method of minimizing $f(\mathbf{y})$ gives in effect a method for solving the linear equation $\mathbf{A}\mathbf{y} = -\mathbf{b}$. The solution of such equations in high dimensions is one of the primary applications of the conjugate gradient method.

We can reduce the problem of minimizing the quadratic function (11.2) to the previous spherical minimization problem by making the change of variables $\mathbf{y} = \mathbf{A}^{-1/2}\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}$ involving the symmetric square root $\mathbf{A}^{-1/2}$ of \mathbf{A}^{-1} . If $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^*$ is the spectral decomposition of the positive definite matrix \mathbf{A} , then $\mathbf{A}^{-1/2} = \mathbf{U}\mathbf{D}^{-1/2}\mathbf{U}^*$. The invertible transformation $\mathbf{x} \mapsto \mathbf{y}$ sends lines into lines and planes into planes. It also sends the function $f(\mathbf{y})$ into the function

$$\begin{aligned} g(\mathbf{x}) &= f(\mathbf{y}) \\ &= f(\mathbf{A}^{-1/2}\mathbf{x} - \mathbf{A}^{-1}\mathbf{b}) \\ &= \frac{1}{2}\|\mathbf{x}\|^2 - \frac{1}{2}\mathbf{b}^*\mathbf{A}^{-1}\mathbf{b} + c \end{aligned}$$

and puts us back where we started, minimizing $\frac{1}{2}\|\mathbf{x}\|^2$. If we have a set of nontrivial orthogonal vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ in \mathbf{x} space, then we can search along each of these directions in turn and achieve the global minima of $g(\mathbf{x})$ and $f(\mathbf{y})$ in n iterations.

For later use, it is important to identify the analogs of the orthogonality conditions (11.1). The direction \mathbf{u}_i in \mathbf{x} space corresponds to the direction \mathbf{v}_i in \mathbf{y} space defined by $\mathbf{A}^{1/2}\mathbf{v}_i = \mathbf{u}_i$. Thus, the condition $\mathbf{u}_i^*\mathbf{u}_j = 0$ for all $j < i$ translates into the condition

$$\mathbf{v}_i^*\mathbf{A}^{1/2}\mathbf{A}^{1/2}\mathbf{v}_j = \mathbf{v}_i^*\mathbf{A}\mathbf{v}_j = 0$$

for all $j < i$. Two vectors \mathbf{v}_i and \mathbf{v}_j satisfying such an orthogonality relation are said to be conjugate. Conjugacy is equivalent to orthogonality under the nonstandard inner product $\mathbf{v}_i^*\mathbf{A}\mathbf{v}_j$. A finite set of conjugate vectors is necessarily linearly independent.

In view of the chain rule, we have $dg(\mathbf{x}_i) = df(\mathbf{y}_i)\mathbf{A}^{-1/2}$ for the point $\mathbf{y}_i = \mathbf{A}^{-1/2}\mathbf{x}_i - \mathbf{A}^{-1}\mathbf{b}$. Thus, the condition $dg(\mathbf{x}_i)\mathbf{u}_j = 0$ for all $j < i$ translates into the condition

$$df(\mathbf{y}_i)\mathbf{A}^{-1/2}\mathbf{A}^{1/2}\mathbf{v}_j = df(\mathbf{y}_i)\mathbf{v}_j = 0 \tag{11.3}$$

for all $j < i$. Alternatively, the condition $df(\mathbf{y}_i)\mathbf{v}_j = 0$ for all $j < i$ is an immediate consequence of the fact that \mathbf{y}_i minimizes the function

$$h(t_1, \dots, t_{i-1}) = f\left(\mathbf{y}_1 + \sum_{j=1}^{i-1} t_j\mathbf{v}_j\right) \tag{11.4}$$

defined on the plane $\mathbf{y}_1 + t_1\mathbf{v}_1 + \dots + t_{i-1}\mathbf{v}_{i-1}$ formed from all linear combinations of the first $i - 1$ search directions.

11.3 The Conjugate Gradient Algorithm

The flaw with this analysis is that it omits any description of how the initial point \mathbf{y}_1 and conjugate directions $\mathbf{v}_1, \dots, \mathbf{v}_n$ are chosen. Choice of \mathbf{y}_1 is

more or less arbitrary, depending on the particular problem and relevant external information. The obvious choice $\mathbf{v}_1 = -\nabla f(\mathbf{y}_1)$ is consistent with an initial search along the direction of steepest descent. At iteration $i > 1$ the conjugate gradient algorithm inductively chooses the search direction

$$\mathbf{v}_i = -\nabla f(\mathbf{y}_i) + \alpha_i \mathbf{v}_{i-1}, \quad (11.5)$$

where

$$\alpha_i = \frac{df(\mathbf{y}_i) \mathbf{A} \mathbf{v}_{i-1}}{\mathbf{v}_{i-1}^* \mathbf{A} \mathbf{v}_{i-1}} \quad (11.6)$$

is defined so that $\mathbf{v}_i^* \mathbf{A} \mathbf{v}_{i-1} = 0$. For $1 \leq j < i - 1$, the conjugacy condition $\mathbf{v}_i^* \mathbf{A} \mathbf{v}_j = 0$ requires

$$0 = -df(\mathbf{y}_i) \mathbf{A} \mathbf{v}_j + \alpha_i \mathbf{v}_{i-1}^* \mathbf{A} \mathbf{v}_j = -df(\mathbf{y}_i) \mathbf{A} \mathbf{v}_j. \quad (11.7)$$

Equality (11.7) is hardly obvious, but we can attack it by noting that in view of definition (11.5) the vectors $\mathbf{v}_1, \dots, \mathbf{v}_{i-1}$ and $\nabla f(\mathbf{y}_1), \dots, \nabla f(\mathbf{y}_{i-1})$ span the same vector subspace. Hence, the condition $df(\mathbf{y}_i) \mathbf{v}_j = 0$ for all $j < i$ is equivalent to the condition $df(\mathbf{y}_i) \nabla f(\mathbf{y}_j) = 0$ for all $j < i$. However, $df(\mathbf{y}_i) \mathbf{v}_j = 0$ for all $j < i$ because \mathbf{y}_i minimizes the function (11.4). Since $\nabla f(\mathbf{y}) = \mathbf{A} \mathbf{y} + \mathbf{b}$ and $\mathbf{y}_{j+1} = \mathbf{y}_j + t_j \mathbf{v}_j$ for some optimal constant t_j , we can also write

$$\nabla f(\mathbf{y}_{j+1}) = \nabla f(\mathbf{y}_j) + t_j \mathbf{A} \mathbf{v}_j. \quad (11.8)$$

It follows that

$$df(\mathbf{y}_i) \mathbf{A} \mathbf{v}_j = \frac{1}{t_j} df(\mathbf{y}_i) [\nabla f(\mathbf{y}_{j+1}) - \nabla f(\mathbf{y}_j)] = 0$$

for $1 \leq j < i - 1$. Except for the details addressed in the next paragraph, this demonstrates equality (11.7) and completes the proof that the search directions $\mathbf{v}_1, \dots, \mathbf{v}_n$ are conjugate.

If at any iteration we have $\nabla f(\mathbf{y}_i) = \mathbf{0}$, then the algorithm terminates with the global minimum. Otherwise, equations (11.3) and (11.5) show that all search directions \mathbf{v}_i satisfy

$$\begin{aligned} df(\mathbf{y}_i) \mathbf{v}_i &= -\|\nabla f(\mathbf{y}_i)\|^2 + \alpha_i df(\mathbf{y}_i) \mathbf{v}_{i-1} \\ &= -\|\nabla f(\mathbf{y}_i)\|^2 \\ &< 0. \end{aligned}$$

As a consequence, $\mathbf{v}_i \neq \mathbf{0}$, $\mathbf{v}_i^* \mathbf{A} \mathbf{v}_i > 0$, and α_{i+1} is well defined. Finally, the inequality $df(\mathbf{y}_i) \mathbf{v}_i < 0$ implies that the search direction \mathbf{v}_i leads downhill from \mathbf{y}_i and that the search constant $t_i > 0$. In fact, the stationarity condition $0 = df(\mathbf{y}_{i+1}) \mathbf{v}_i$ and equation (11.8) lead to the conclusion

$$t_i = -\frac{df(\mathbf{y}_i) \mathbf{v}_i}{\mathbf{v}_i^* \mathbf{A} \mathbf{v}_i} = -\frac{(\mathbf{A} \mathbf{y}_i + \mathbf{b})^* \mathbf{v}_i}{\mathbf{v}_i^* \mathbf{A} \mathbf{v}_i}. \quad (11.9)$$

In generalizing the conjugate gradient algorithm to non-quadratic functions, we preserve most of the structure of the algorithm. Thus, the revised algorithm first searches along the negative gradient $\mathbf{v}_1 = -\nabla f(\mathbf{y}_1)$ emanating from the initial point \mathbf{y}_1 . At iteration i it searches along the direction \mathbf{v}_i defined by equality (11.5), avoiding the explicit formula (11.9) for t_i . The formula (11.6) for α_i is problematic because it appears to require \mathbf{A} . Hestenes and Stiefel recommend the alternative formula

$$\alpha_i = \frac{df(\mathbf{y}_i)[\nabla f(\mathbf{y}_i) - \nabla f(\mathbf{y}_{i-1})]}{\mathbf{v}_{i-1}^*[\nabla f(\mathbf{y}_i) - \nabla f(\mathbf{y}_{i-1})]} \quad (11.10)$$

based on the substitution

$$\mathbf{A}\mathbf{v}_{i-1} = \frac{1}{t_{i-1}}[\nabla f(\mathbf{y}_i) - \nabla f(\mathbf{y}_{i-1})].$$

Polak and Ribière suggest the further substitutions

$$\begin{aligned} \mathbf{v}_{i-1}^* \nabla f(\mathbf{y}_i) &= 0 \\ \mathbf{v}_{i-1}^* \nabla f(\mathbf{y}_{i-1}) &= [-df(\mathbf{y}_{i-1}) + \alpha_{i-1} \mathbf{v}_{i-2}] \nabla f(\mathbf{y}_{i-1}) \\ &= -\|\nabla f(\mathbf{y}_{i-1})\|^2. \end{aligned}$$

These produce the second alternative

$$\alpha_i = \frac{df(\mathbf{y}_i)[\nabla f(\mathbf{y}_i) - \nabla f(\mathbf{y}_{i-1})]}{\|\nabla f(\mathbf{y}_{i-1})\|^2}. \quad (11.11)$$

Finally, Fletcher and Reeves note that the identity

$$\begin{aligned} df(\mathbf{y}_i) \nabla f(\mathbf{y}_{i-1}) &= df(\mathbf{y}_i)[- \mathbf{v}_{i-1} + \alpha_{i-1} \mathbf{v}_{i-2}] \\ &= 0 \end{aligned}$$

yields the third alternative

$$\alpha_i = \frac{\|\nabla f(\mathbf{y}_i)\|^2}{\|\nabla f(\mathbf{y}_{i-1})\|^2}. \quad (11.12)$$

Almost no one now uses the Hestenes-Stiefel update (11.10). Current opinion is divided between the Polak-Ribière update (11.11) and the Fletcher-Reeves update (11.12). *Numerical Recipes* [215] codes both formulas but leans toward the Polak-Ribière formula. In addition to this issue, there are other practical concerns in implementing the conjugate gradient algorithm. For example, if we fail to stop once the gradient $\nabla f(\mathbf{y})$ vanishes, then the Polak-Ribière and Fletcher-Reeves updates are undefined. This suggests stopping when $\|\nabla f(\mathbf{y}_i)\| \leq \epsilon \|\nabla f(\mathbf{y}_1)\|$ for some small $\epsilon > 0$. There is also the problem of loss of conjugacy. Assuming $f(\mathbf{y})$ is defined on \mathbb{R}^n , it is common practice to restart the conjugate gradient algorithm with

the steepest descent direction every n iterations. This is also a good idea whenever the descent condition $df(\mathbf{y}_i)\mathbf{v}_i < 0$ fails. Finally, the algorithm is incomplete without specifying a line search algorithm. The next section discusses some of the ways of conducting a line search. The references [2, 92, 215] provide a fuller account and appropriate computer code.

11.4 Line Search Methods

The secant method of Example 10.2.4 can obviously be adapted to minimize the objective function $f(\mathbf{y})$ along a line $r \mapsto \mathbf{x} + r\mathbf{v}$ in \mathbb{R}^n . If we define $g(r) = f(\mathbf{x} + r\mathbf{v})$, then we proceed by searching for a zero of the derivative $g'(r) = df(\mathbf{x} + r\mathbf{v})\mathbf{v}$. In this guise, the secant method is known as the method of false position. It iterates according to

$$r_{m+1} = r_m - \frac{g'(r_m)(r_{m-1} - r_m)}{g'(r_{m-1}) - g'(r_m)}.$$

Two criticisms of the method of false position immediately come to mind. One is that it indiscriminately heads for maxima as well as minima. Another is that it does not make full use of the available information.

A better alternative to the method of false position is to approximate $g(r)$ by a cubic polynomial matching the values of $g(r)$ and $g'(r)$ at r_m and r_{m-1} . Minimizing the cubic should lead to an improved estimate r_{m+1} of the minimum of $g(r)$. It simplifies matters notationally to rescale the interval by setting $h(s) = g(r_{m-1} + sd_m)$ and $d_m = r_m - r_{m-1}$. Now $s = 0$ corresponds to r_{m-1} and $s = 1$ corresponds to r_m . Furthermore, the chain rule implies $h'(s) = g'(r_{m-1} + sd_m)d_m$. Given these conventions, the theory of Hermite interpolation [123] suggests approximating $h(s)$ by the cubic polynomial

$$\begin{aligned} p(s) &= (s-1)^2 h_0 + s^2 h_1 + s(s-1)[(s-1)(h'_0 + 2h_0) + s(h'_1 - 2h_1)] \\ &= (2h_0 + h'_0 - 2h_1 + h'_1)s^3 + (-3h_0 - 2h'_0 + 3h_1 - h'_1)s^2 + h'_0 s + h_0, \end{aligned}$$

where $h_0 = h(0)$, $h'_0 = h'(0)$, $h_1 = h(1)$, and $h'_1 = h'(1)$. One can readily verify that $p(0) = h_0$, $p'(0) = h'_0$, $p(1) = h_1$, and $p'(1) = h'_1$.

The conjugate gradient method is locally descending in the sense that $p'(0) = h'_0 < 0$. To be on the cautious side, $p'(1) = h'_1 > 0$ should hold and $p(s)$ should be convex throughout the interval $[0, 1]$. To check convexity, it suffices to check the conditions $p''(0) \geq 0$ and $p''(1) \geq 0$ since $p''(s)$ is linear. Straightforward calculation shows that

$$\begin{aligned} p''(0) &= -6h_0 + 6h_1 - 4h'_0 - 2h'_1 \\ p''(1) &= 6h_0 - 6h_1 + 2h'_0 + 4h'_1. \end{aligned}$$

Thus, $p(s)$ is convex throughout $[0, 1]$ if and only if

$$\frac{1}{3}h'_1 + \frac{2}{3}h'_0 \leq h_1 - h_0 \leq \frac{2}{3}h'_1 + \frac{1}{3}h'_0. \quad (11.13)$$

Under these conditions, a local minimum of $p(s)$ occurs on $[0, 1]$. The pertinent root of the two possible roots of $p'(s) = 0$ is determined by the sign of the coefficient $2h_0 + h'_0 - 2h_1 + h'_1$ of s^3 in $p(s)$. If this coefficient is positive, then the right root furnishes the minimum, and if this coefficient is negative, then the left root furnishes the minimum. The two roots can be calculated simultaneously by solving the quadratic equation

$$\begin{aligned} p'(s) &= 3(2h_0 + h'_0 - 2h_1 + h'_1)s^2 + 2(-3h_0 - 2h'_0 + 3h_1 - h'_1)s + h'_0 \\ &= 0. \end{aligned}$$

If the condition $p'(1) = h'_1 > 0$ or the convexity conditions (11.13) fail, or if the minimum of the cubic leads to an increase in $g(r)$, then one should fall back on more conservative search methods. Golden search involves recursively bracketing a minimum by three points $a < b < c$ satisfying $g(b) < \min\{g(a), g(c)\}$. The analogous method of bisection brackets a zero of $g(r)$ by two points $a < b$ satisfying $g(a)g(b) < 0$. For the moment we ignore the question of how the initial three points a , b , and c are chosen in golden search.

To replace the bracketing interval (a, c) by a shorter interval, we choose $d \in (a, c)$ so that d belongs to the longer of the two intervals (a, b) and (b, c) . Without loss of generality, suppose $b < d < c$. If $g(d) < g(b)$, then the three points $b < d < c$ bracket a minimum. If $g(d) > g(b)$, then the three points $a < b < d$ bracket a minimum. In the case of a tie $g(d) = g(b)$, we choose $b < d < c$ when $g(c) < g(a)$ and $a < b < d$ when $g(a) < g(c)$.

These sensible rules do not address the problem of choosing d . Consider the fractional distances

$$\beta = \frac{b-a}{c-a}, \quad \delta = \frac{d-b}{c-a}$$

along the interval (a, c) . The next bracketing interval will have a fractional length of either $1 - \beta$ or $\beta + \delta$. To guard against the worst case, we should take $1 - \beta = \beta + \delta$. This determines $\delta = 1 - 2\beta$ and hence d . One could leave matters as they now stand, but the argument is taken one step farther in golden search. If we imagine repeatedly performing golden search, then scale similarity is expected to set in eventually so that

$$\beta = \frac{b-a}{c-a} = \frac{d-b}{c-b} = \frac{\delta}{1-\beta}.$$

Substituting $\delta = 1 - 2\beta$ in this identity and cross multiplying give the quadratic $\beta^2 - 3\beta + 1 = 0$ with solution

$$\beta = \frac{3 - \sqrt{5}}{2}$$

equal to the golden mean of ancient Greek mathematics. Following this reasoning, we should take $\delta = \sqrt{5} - 2 = 0.2361$.

There is little theory to guide us in finding an initial bracketing triple $a < b < c$. It is clear that $a = 0$ is one natural choice. In view of the condition $g'(0) < 0$, the point b can be chosen close to 0 as well. This leaves c , which is usually selected based on specific knowledge of $f(y)$, parabolic extrapolation, or repeated doubling of some small arbitrary distance.

11.5 Stopping Criteria

Deciding when to terminate an iterative method is more subtle than it might seem. In solving a one-dimensional nonlinear equation $g(x) = 0$, there are basically two tests. One can declare convergence when $|g(x_n)|$ is small or when x_n does not change much from one iteration to the next. Ideally, both tests should be satisfied. However, there are questions of scale. Our notion of small depends on the typical magnitudes of $g(x)$ and x , and stopping criteria should reflect these magnitudes [66]. Suppose $a > 0$ represents the typical magnitude of $g(x)$. Then a sensible criterion of the first kind is to stop when $|g(x_n)| < \epsilon a$ for $\epsilon > 0$ small. If $b > 0$ represents the typical magnitude of x , then a sensible criterion of the second kind is to stop when

$$|x_n - x_{n-1}| \leq \epsilon \max\{|x_n|, b\}. \quad (11.14)$$

To achieve p significant digits in the solution x_∞ , take $\epsilon = 10^{-p}$.

When we optimize a function $f(x)$ with derivative $g(x) = f'(x)$, a third test comes into play. Now it is desirable for $f(x)$ to remain relatively constant near convergence. If $c > 0$ represents the typical magnitude of $f(x)$, then our final stopping criterion is

$$|f(x_n) - f(x_{n-1})| \leq \epsilon \max\{|f(x_n)|, c\}.$$

The second and third criteria generalize better than the first criterion to higher-dimensional problems because solutions often occur on boundaries or manifolds where the gradient $\nabla f(\mathbf{x})$ is not required to vanish. The Karush-Kuhn-Tucker conditions are acceptable substitute for Fermat's condition provided the Lagrange multipliers are known. In higher dimensions, one should apply the criterion (11.14) to each coordinate of \mathbf{x} . Choice of the typical magnitudes a , b , and c is problem specific, and some optimization programs leave this up to the discretion of the user. Often problems can be rescaled by an appropriate choice of units so that the choice $a = b = c = 1$ is reasonable. When in doubt about typical magnitudes, take this default and check whether the output of a preliminary computer run justifies the assumption.

11.6 Quasi-Newton Methods

Quasi-Newton methods of minimization update the current approximation \mathbf{H}_i to the second differential $d^2f(\mathbf{x}_i)$ of the objective function $f(\mathbf{x})$ by a low-rank perturbation satisfying a secant condition. The secant condition originates from the first-order Taylor approximation

$$\nabla f(\mathbf{x}_{i+1}) - \nabla f(\mathbf{x}_i) \approx d^2f(\mathbf{x}_i)(\mathbf{x}_{i+1} - \mathbf{x}_i).$$

If we set

$$\begin{aligned} \mathbf{g}_i &= \nabla f(\mathbf{x}_{i+1}) - \nabla f(\mathbf{x}_i) \\ \mathbf{d}_i &= \mathbf{x}_{i+1} - \mathbf{x}_i, \end{aligned}$$

then the secant condition reads $\mathbf{H}_{i+1}\mathbf{d}_i = \mathbf{g}_i$. The unique, symmetric, rank-one update to \mathbf{H}_i satisfying the secant condition is furnished by Davidon's formula [56]

$$\mathbf{H}_{i+1} = \mathbf{H}_i + c_i \mathbf{v}_i \mathbf{v}_i^* \quad (11.15)$$

with the constant c_i and the vector \mathbf{v}_i specified by

$$\begin{aligned} c_i &= -\frac{1}{(\mathbf{H}_i \mathbf{d}_i - \mathbf{g}_i)^* \mathbf{d}_i} \\ \mathbf{v}_i &= \mathbf{H}_i \mathbf{d}_i - \mathbf{g}_i. \end{aligned} \quad (11.16)$$

An immediate concern is that the constant c_i is undefined when the inner product $(\mathbf{H}_i \mathbf{d}_i - \mathbf{g}_i)^* \mathbf{d}_i = 0$. In such situations or when

$$|(\mathbf{H}_i \mathbf{d}_i - \mathbf{g}_i)^* \mathbf{d}_i| \ll \|\mathbf{H}_i \mathbf{d}_i - \mathbf{g}_i\| \cdot \|\mathbf{d}_i\|,$$

then the secant adjustment is ignored, and the value \mathbf{H}_i is retained for \mathbf{H}_{i+1} .

We have stressed the desirability of maintaining a positive definite approximation \mathbf{H}_i to the second differential $d^2f(\mathbf{x}_i)$. Because this is not always possible with the rank-one update, numerical analysts have investigated rank-two updates. The involvement of the vectors \mathbf{g}_i and $\mathbf{H}_i \mathbf{d}_i$ in the rank-one update suggests trying a rank-two update of the form

$$\mathbf{H}_{i+1} = \mathbf{H}_i + b_i \mathbf{g}_i \mathbf{g}_i^* + c_i \mathbf{H}_i \mathbf{d}_i \mathbf{d}_i^* \mathbf{H}_i. \quad (11.17)$$

Taking the product of both sides of this equation with \mathbf{d}_i gives

$$\mathbf{H}_{i+1} \mathbf{d}_i = \mathbf{H}_i \mathbf{d}_i + b_i \mathbf{g}_i \mathbf{g}_i^* \mathbf{d}_i + c_i \mathbf{H}_i \mathbf{d}_i \mathbf{d}_i^* \mathbf{H}_i \mathbf{d}_i.$$

To achieve consistency with the secant condition $\mathbf{H}_{i+1} \mathbf{d}_i = \mathbf{g}_i$, we set

$$b_i = \frac{1}{\mathbf{g}_i^* \mathbf{d}_i}, \quad c_i = -\frac{1}{\mathbf{d}_i^* \mathbf{H}_i \mathbf{d}_i}.$$

The resulting rank-two update was proposed by Broyden, Fletcher, Goldfarb, and Shanno and is consequently known as the BFGS update.

The symmetric rank-one update (11.15) certainly preserves positive definiteness when $c_i \geq 0$. If $c_i < 0$, then \mathbf{H}_{i+1} is positive definite only if

$$\begin{aligned} \mathbf{v}_i^* \mathbf{H}_i^{-1} [\mathbf{H}_i + c_i \mathbf{v}_i \mathbf{v}_i^*] \mathbf{H}_i^{-1} \mathbf{v}_i &= \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i [1 + c_i \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i] \\ &> 0. \end{aligned}$$

In other words, the condition

$$1 + c_i \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i > 0 \quad (11.18)$$

must hold. Conversely, condition (11.18) is sufficient to guarantee positive definiteness of \mathbf{H}_{i+1} . This fact can be most easily demonstrated by noting the Sherman-Morrison inversion formula [197]

$$[\mathbf{H}_i + c_i \mathbf{v}_i \mathbf{v}_i^*]^{-1} = \mathbf{H}_i^{-1} - \frac{c_i}{1 + c_i \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i} \mathbf{H}_i^{-1} \mathbf{v}_i [\mathbf{H}_i^{-1} \mathbf{v}_i]^*. \quad (11.19)$$

Formula (11.19) shows that $[\mathbf{H}_i + c_i \mathbf{v}_i \mathbf{v}_i^*]^{-1}$ exists and is positive definite under condition (11.18). Since the inverse of a positive definite matrix is positive definite, it follows that $\mathbf{H}_i + c_i \mathbf{v}_i \mathbf{v}_i^*$ is positive definite as well.

If $c_i < 0$ in the rank-one update but condition (11.18) fails, then there are various options. The preferred is implementation of the trust region strategy discussed in Sect. 11.7. Alternatively, one can shrink c_i to maintain positive definiteness. Unfortunately, condition (11.18) gives too little guidance. Problem 12 shows how to control the size of $\det \mathbf{H}_{i+1}$ while simultaneously forcing positive definiteness. An even better strategy that monitors the condition number of \mathbf{H}_{i+1} rather than $\det \mathbf{H}_{i+1}$ is sketched in Problem 14. Finally, there is the option of using c_i as defined but perturbing \mathbf{H}_{i+1} by adding a constant multiple $\mu \mathbf{I}$ of the identity matrix. This tactic is similar in spirit to the trust region method. If λ_1 is the smallest eigenvalue of \mathbf{H}_{i+1} , then $\mathbf{H}_{i+1} + \mu \mathbf{I}$ is positive definite whenever $\lambda_1 + \mu > 0$. Problem 15 discusses a fast algorithm for finding λ_1 . With appropriate safeguards, some numerical analysts [51, 152] consider the rank-one update superior to the BFGS update.

Positive definiteness is almost automatic with the BFGS update (11.17). The key turns out to be the inequality

$$0 < \mathbf{g}_i^* \mathbf{d}_i = df(\mathbf{x}_{i+1}) \mathbf{d}_i - df(\mathbf{x}_i) \mathbf{d}_i. \quad (11.20)$$

This is ordinarily true for two reasons. First, because \mathbf{d}_i is proportional to the current search direction $\mathbf{v}_i = -\mathbf{H}_i^{-1} \nabla f(\mathbf{x}_i)$, positive definiteness of \mathbf{H}_i^{-1} implies $-df(\mathbf{x}_i) \mathbf{d}_i > 0$. Second, when a full search is conducted, the identity $df(\mathbf{x}_{i+1}) \mathbf{d}_i = 0$ holds. Even a partial search typically entails condition (11.20). Section 12.6 takes up the issue of partial line searches.

The speed of partial line searches compared to that of full line searches makes quasi-Newton methods superior to the conjugate gradient method on small-scale problems.

To show that the BFGS update \mathbf{H}_{i+1} is positive definite when condition (11.20) holds, we examine the quadratic form

$$\mathbf{u}^* \mathbf{H}_{i+1} \mathbf{u} = \mathbf{u}^* \mathbf{H}_i \mathbf{u} + \frac{(\mathbf{g}_i^* \mathbf{u})^2}{\mathbf{g}_i^* \mathbf{d}_i} - \frac{(\mathbf{u}^* \mathbf{H}_i \mathbf{d}_i)^2}{\mathbf{d}_i^* \mathbf{H}_i \mathbf{d}_i} \quad (11.21)$$

for $\mathbf{u} \neq \mathbf{0}$. Applying Cauchy's inequality to the vectors $\mathbf{a} = \mathbf{H}_i^{1/2} \mathbf{u}$ and $\mathbf{b} = \mathbf{H}_i^{1/2} \mathbf{d}_i$ gives

$$(\mathbf{u}^* \mathbf{H}_i \mathbf{d}_i)^2 \leq (\mathbf{u}^* \mathbf{H}_i \mathbf{u})(\mathbf{d}_i^* \mathbf{H}_i \mathbf{d}_i),$$

with equality if and only if \mathbf{u} is proportional to \mathbf{d}_i . Hence, the sum of the first and third terms on the right of equality (11.21) is nonnegative. In the event that \mathbf{u} is proportional to \mathbf{d}_i , the second term on the right of equality (11.21) is positive by assumption. It follows that $\mathbf{u}^* \mathbf{H}_{i+1} \mathbf{u} > 0$ and therefore that \mathbf{H}_{i+1} is positive definite.

In successful applications of quasi-Newton methods, choice of the initial matrix \mathbf{H}_1 is critical. Setting $\mathbf{H}_1 = \mathbf{I}$ is convenient but often poorly scaled for a particular problem. In maximum likelihood estimation, the expected information matrix $J(\mathbf{x}_1)$, if available, is preferable to the identity matrix. In some problems, $J(\mathbf{x})$ is cheap to compute and manipulate for special values of \mathbf{x} . For instance, $J(\mathbf{x})$ may be diagonal in certain circumstances. These special \mathbf{x} should be considered as starting points for a quasi-Newton search.

It is possible to carry forward approximations \mathbf{K}_i of $d^2 f(\mathbf{x}_i)^{-1}$ rather than of $d^2 f(\mathbf{x}_i)$. This tactic has the advantage of avoiding matrix inversion in computing the quasi-Newton search direction $\mathbf{v}_i = -\mathbf{K}_i \nabla f(\mathbf{x}_i)$. The basic idea is to restate the secant condition $\mathbf{H}_{i+1} \mathbf{d}_i = \mathbf{g}_i$ as the inverse secant condition $\mathbf{K}_{i+1} \mathbf{g}_i = \mathbf{d}_i$. This substitution leads to the symmetric rank-one update

$$\mathbf{K}_{i+1} = \mathbf{K}_i + c_i \mathbf{w}_i \mathbf{w}_i^*, \quad (11.22)$$

where $c_i = -[(\mathbf{K}_i \mathbf{g}_i - \mathbf{d}_i)^* \mathbf{g}_i]^{-1}$ and $\mathbf{w}_i = \mathbf{K}_i \mathbf{g}_i - \mathbf{d}_i$. Note that monitoring positive definiteness of \mathbf{K}_i is still an issue.

For a rank-two update, our earlier arguments apply provided we interchange the roles of \mathbf{d}_i and \mathbf{g}_i . The Davidon-Fletcher-Powell (DFP) update

$$\mathbf{K}_{i+1} = \mathbf{K}_i + b_i \mathbf{d}_i \mathbf{d}_i^* + c_i \mathbf{K}_i \mathbf{g}_i \mathbf{g}_i^* \mathbf{K}_i \quad (11.23)$$

with

$$b_i = \frac{1}{\mathbf{g}_i^* \mathbf{d}_i}, \quad c_i = -\frac{1}{\mathbf{g}_i^* \mathbf{K}_i \mathbf{g}_i}$$

is a competitor to the BFGS update, but the consensus seems to be that the BFGS update is superior to the DFP update in practice [66].

In closing this section, we would like to prove that the BFGS algorithm with an exact line search converges in n or fewer iterations for the strictly convex quadratic function (11.2) defined on \mathbb{R}^n . Recall that at iteration i we search along the direction $\mathbf{v}_i = -\mathbf{H}_i^{-1}\nabla f(\mathbf{x}_i)$ and then in preparation for the next iteration construct \mathbf{H}_{i+1} according to the BFGS formula (11.17). Unless $\nabla f(\mathbf{x}_i) = \mathbf{0}$ and the iterates converge prematurely, the current increment \mathbf{d}_i is a positive multiple $t_i\mathbf{v}_i$ of the search direction \mathbf{v}_i . Our proof of convergence consists of a subtle inductive argument proving three claims in parallel. These claims amount to the conjugacy condition $\mathbf{v}_{i+1}\mathbf{A}\mathbf{v}_j = 0$, the extended secant condition $\mathbf{H}_{i+1}\mathbf{d}_j = \mathbf{g}_j$, and the gradient perpendicularity condition $df(\mathbf{x}_{i+1})\mathbf{v}_j = 0$, each for all $1 \leq j \leq i$ and all $i \leq n$. Given the efficacy of successive searches along conjugate directions as demonstrated in Sect. 11.2, the conjugacy condition $\mathbf{v}_{i+1}\mathbf{A}\mathbf{v}_j = 0$ guarantees convergence to the minimum of $f(\mathbf{y})$ in n or fewer iterations.

The case $i = 0$, where all three conditions are vacuous, gets the induction on i started. In general, assume that the three conditions are true for $i - 1$ and all $1 \leq j \leq i - 1$. Equation (11.3) validates the gradient perpendicularity condition for any set of conjugate directions $\mathbf{v}_1, \dots, \mathbf{v}_i$, not just the ones determined by the BFGS update (11.17). Given the gradient identity $\mathbf{A}\mathbf{d}_j = \mathbf{g}_j$ and the validity of the extended secant condition $\mathbf{H}_{i+1}\mathbf{d}_j = \mathbf{g}_j$, we calculate

$$\begin{aligned} \mathbf{v}_{i+1}^*\mathbf{A}\mathbf{v}_j &= -df(\mathbf{x}_{i+1})\mathbf{H}_{i+1}^{-1}\mathbf{A}\mathbf{v}_j \\ &= -t_j^{-1}df(\mathbf{x}_{i+1})\mathbf{H}_{i+1}^{-1}\mathbf{A}\mathbf{d}_j \\ &= -t_j^{-1}df(\mathbf{x}_{i+1})\mathbf{H}_{i+1}^{-1}\mathbf{g}_j \\ &= -t_j^{-1}df(\mathbf{x}_{i+1})\mathbf{d}_j \\ &= -df(\mathbf{x}_{i+1})\mathbf{v}_j \\ &= 0, \end{aligned} \tag{11.24}$$

which is the required conjugacy condition.

Thus, it suffices to prove $\mathbf{H}_{i+1}\mathbf{d}_j = \mathbf{g}_j$ for $j \leq i$. The case $j = i$ is just the ordinary secant requirement. For $j < i$, we observe that

$$\mathbf{H}_{i+1}\mathbf{d}_j = \mathbf{H}_i\mathbf{d}_j + b_i\mathbf{g}_i\mathbf{g}_i^*\mathbf{d}_j + c_i\mathbf{H}_i\mathbf{d}_i\mathbf{d}_i^*\mathbf{H}_i\mathbf{d}_j.$$

Now the equalities $\mathbf{H}_i\mathbf{d}_j = \mathbf{g}_j$ and $\mathbf{g}_i^*\mathbf{d}_j = \mathbf{d}_i^*\mathbf{A}\mathbf{d}_j = 0$ hold by the induction hypothesis. Likewise,

$$\begin{aligned} \mathbf{d}_i^*\mathbf{H}_i\mathbf{d}_j &= \mathbf{d}_i^*\mathbf{g}_j \\ &= \mathbf{d}_i^*\mathbf{A}\mathbf{d}_j \\ &= 0 \end{aligned}$$

follows from the induction hypothesis. Combining these equalities makes it clear that $\mathbf{H}_{i+1}\mathbf{d}_j = \mathbf{g}_j$ and completes the induction and the proof.

11.7 Trust Regions

If the quadratic approximation

$$f(\mathbf{x}) \approx f(\mathbf{x}_i) + df(\mathbf{x}_i)(\mathbf{x} - \mathbf{x}_i) + \frac{1}{2}(\mathbf{x} - \mathbf{x}_i)^* \mathbf{H}_i(\mathbf{x} - \mathbf{x}_i)$$

to the objective function $f(\mathbf{x})$ is poor, then the naive step

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \mathbf{H}_i^{-1} \nabla f(\mathbf{x}_i) \quad (11.25)$$

computed in a quasi-Newton method may be absurdly large. This situation often occurs for early iterates. One remedy is to minimize the quadratic approximation to $f(\mathbf{x})$ subject to the spherical constraint $\|\mathbf{x} - \mathbf{x}_i\|^2 \leq r^2$ for a fixed radius r . This constrained optimization problem has a solution regardless of whether \mathbf{H}_i is positive definite, but to simplify matters in the remaining discussion, we assume that \mathbf{H}_i is positive definite. According to Proposition 5.2.1, the solution satisfies the multiplier rule

$$\mathbf{0} = \nabla f(\mathbf{x}_i) + \mathbf{H}_i(\mathbf{x} - \mathbf{x}_i) + \mu(\mathbf{x} - \mathbf{x}_i) \quad (11.26)$$

for some nonnegative constant μ . If the point \mathbf{x}_{i+1} generated by the ordinary step (11.25) occurs within the open ball $\|\mathbf{x} - \mathbf{x}_i\| < r$, then the multiplier $\mu = 0$. Of course, there is no guarantee that this choice of \mathbf{x}_{i+1} will lead to a decrease in $f(\mathbf{x})$. If it does not, then one should reduce r , for instance to $\frac{r}{2}$, and try again.

When the point \mathbf{x}_{i+1} generated by the ordinary step (11.25) occurs outside the open ball $\|\mathbf{x} - \mathbf{x}_i\| < r$, we are obliged to look for a minimum point \mathbf{x} to the quadratic approximation satisfying $\|\mathbf{x} - \mathbf{x}_i\| = r$. In this case the multiplier μ may be positive. The fact that μ is unknown makes it impossible to find the minimum in closed form. In principle, one can overcome this difficulty by solving for μ iteratively, say by Newton's method. Hence, we view equation (11.26) as defining $\mathbf{x} - \mathbf{x}_i$ as a function of μ and ask for the value of μ that yields $\|\mathbf{x} - \mathbf{x}_i\| = r$. To simplify notation, let $\mathbf{H} = \mathbf{H}_i$, $\mathbf{y} = \mathbf{x} - \mathbf{x}_i$, and $\mathbf{e} = -\nabla f(\mathbf{x}_i)$. We now seek a zero of the function

$$\phi(\mu) = \frac{1}{r} - \frac{1}{\|\mathbf{y}(\mu)\|} \quad (11.27)$$

with $\mathbf{y}(\mu)$ defined by $(\mathbf{H} + \mu\mathbf{I})\mathbf{y}(\mu) = \mathbf{e}$. Note that $\phi(0) > 0$ if and only if $\|\mathbf{y}(0)\| > r$. To implement Newton's method, we need $\phi'(\mu)$. An easy calculation shows that

$$\phi'(\mu) = \frac{\mathbf{y}(\mu)^* \mathbf{y}'(\mu)}{\|\mathbf{y}(\mu)\|^3}.$$

Unfortunately, this formula contains the unknown derivative $\mathbf{y}'(\mu)$. However, differentiation of the equation $(\mathbf{H} + \mu\mathbf{I})\mathbf{y}(\mu) = \mathbf{e}$ readily yields

$$\mathbf{y}(\mu) + (\mathbf{H} + \mu\mathbf{I})\mathbf{y}'(\mu) = \mathbf{0}, \quad (11.28)$$

which implies $\mathbf{y}'(\mu) = -(\mathbf{H} + \mu\mathbf{I})^{-1}\mathbf{y}(\mu)$. The complete formula

$$\phi'(\mu) = -\frac{\mathbf{y}(\mu)^*(\mathbf{H} + \mu\mathbf{I})^{-1}\mathbf{y}(\mu)}{\|\mathbf{y}(\mu)\|^3}$$

shows that $\phi(\mu)$ is strictly decreasing. Problem 16 asks the reader to calculate $\phi''(\mu)$ and verify that it is nonnegative. Problem 17 asserts that $\phi(\mu)$ is negative for large μ . Hence, there is a unique Lagrange multiplier $\mu_i > 0$ solving $\phi(\mu) = 0$ whenever $\phi(0) > 0$. The corresponding $\mathbf{y}(\mu_i)$ solves the trust region problem [198, 241].

If one is willing to extract the spectral decomposition of $\mathbf{U}\mathbf{D}\mathbf{U}^t$ of \mathbf{H}_i , then the process can be simplified. Let $\mathbf{z} = \mathbf{U}^t(\mathbf{x} - \mathbf{x}_i)$ and $\mathbf{b} = \mathbf{U}^t\nabla f(\mathbf{x}_i)$. Then the trust region problem reduces to minimizing $\frac{1}{2}\mathbf{z}^t\mathbf{D}\mathbf{z} + \mathbf{b}^t\mathbf{z}$ subject to $\|\mathbf{z}\|^2 \leq r^2$. The stationarity conditions for the corresponding Lagrangian

$$\mathcal{L}(\mathbf{z}, \mu) = \frac{1}{2}\mathbf{z}^t\mathbf{D}\mathbf{z} + \mathbf{b}^t\mathbf{z} + \frac{\mu}{2}(\|\mathbf{z}\|^2 - r^2)$$

yield

$$z_j = -\frac{b_j}{d_j + \mu}.$$

where d_j is the j th diagonal entry of \mathbf{D} . When $\mathbf{z} = -\mathbf{D}^{-1}\mathbf{b}$ satisfies the constraint $\|\mathbf{z}\|^2 \leq r^2$, we take $\mu = 0$. Otherwise, we solve the constraint equality

$$r^2 = \sum_j \left(\frac{b_j}{d_j + \mu} \right)^2$$

for μ numerically and determine \mathbf{z} and $\mathbf{x} = \mathbf{U}\mathbf{z} + \mathbf{x}_i$ accordingly. For more details about trust regions and their practical implementation, see the books [66, 107, 151, 205].

11.8 Problems

1. Suppose you possess n conjugate vectors $\mathbf{v}_1, \dots, \mathbf{v}_n$ for the $n \times n$ positive definite matrix \mathbf{A} . Describe how you can use the expansion $\mathbf{x} = \sum_{i=1}^n c_i \mathbf{v}_i$ to solve the linear equation $\mathbf{A}\mathbf{x} = \mathbf{b}$.
2. Suppose that \mathbf{A} is an $n \times n$ positive definite matrix and that the nontrivial vectors $\mathbf{u}_1, \dots, \mathbf{u}_n$ satisfy

$$\mathbf{u}_i^* \mathbf{A} \mathbf{u}_j = 0, \quad \mathbf{u}_i^* \mathbf{u}_j = 0$$

for all $i \neq j$. Demonstrate that the \mathbf{u}_i are eigenvectors of \mathbf{A} .

3. Suppose that the $n \times n$ symmetric matrix \mathbf{A} satisfies $\mathbf{v}^* \mathbf{A} \mathbf{v} \neq 0$ for all $\mathbf{v} \neq \mathbf{0}$ and that $\{\mathbf{u}_1, \dots, \mathbf{u}_n\}$ is a basis of \mathbb{R}^n . If one defines $\mathbf{v}_1 = \mathbf{u}_1$ and inductively

$$\mathbf{v}_k = \mathbf{u}_k - \sum_{j=1}^{k-1} \frac{\mathbf{u}_k^* \mathbf{A} \mathbf{v}_j}{\mathbf{v}_j^* \mathbf{A} \mathbf{v}_j} \mathbf{v}_j$$

for $k = 2, \dots, n$, then show that the vectors $\{\mathbf{v}_1, \dots, \mathbf{v}_n\}$ are conjugate and provide a basis of \mathbb{R}^n . Note that \mathbf{A} need not be positive definite.

4. Consider extracting a root of the equation $x^2 = 0$ by Newton's method starting from $x_0 = 1$. Show that it is impossible to satisfy the convergence criterion

$$|x_n - x_{n-1}| \leq \epsilon \max\{|x_n|, |x_{n-1}|\}$$

for $\epsilon = 10^{-7}$ [66]. This example favors the alternative stopping rule (11.14).

5. Let \mathbf{M} be an $n \times n$ matrix and \mathbf{d} and \mathbf{g} be $n \times 1$ vectors. Show that the matrix

$$\mathbf{N}_{\text{opt}} = \mathbf{M} + \|\mathbf{d}\|^{-2} (\mathbf{g} - \mathbf{M} \mathbf{d}) \mathbf{d}^*$$

minimizes the distance $\|\mathbf{N} - \mathbf{M}\|$ between \mathbf{M} and an arbitrary $n \times n$ matrix \mathbf{N} subject to the secant condition $\mathbf{N} \mathbf{d} = \mathbf{g}$ [66]. Unfortunately, the rank-one update \mathbf{N}_{opt} is not symmetric when \mathbf{M} is symmetric. (Hints: Note that $(\mathbf{N} - \mathbf{M}) \mathbf{d} = (\mathbf{N}_{\text{opt}} - \mathbf{M}) \mathbf{d}$ for every such \mathbf{N} and that the outer product $\mathbf{d} \mathbf{d}^*$ has induced matrix norm $\|\mathbf{d}\|^2$.)

6. Let \mathbf{M} be an $n \times n$ symmetric matrix and \mathbf{d} and \mathbf{g} be $n \times 1$ vectors. Powell proposed the rank-two update

$$\mathbf{N}_{\text{opt}} = \mathbf{M} + \frac{(\mathbf{g} - \mathbf{M} \mathbf{d}) \mathbf{d}^* + \mathbf{d} (\mathbf{g} - \mathbf{M} \mathbf{d})^*}{\|\mathbf{d}\|^2} - \frac{(\mathbf{g} - \mathbf{M} \mathbf{d})^* \mathbf{d} \mathbf{d}^*}{\|\mathbf{d}\|^4}$$

to \mathbf{M} . Show that \mathbf{N}_{opt} is symmetric, has rank two, and satisfies the secant condition $\mathbf{N}_{\text{opt}} \mathbf{d} = \mathbf{g}$ [66].

7. Continuing Problem 6, show that the matrix \mathbf{N}_{opt} minimizes the distance $\|\mathbf{N} - \mathbf{M}\|_F$ between \mathbf{M} and an arbitrary $n \times n$ symmetric matrix \mathbf{N} subject to the secant condition $\mathbf{N} \mathbf{d} = \mathbf{g}$. Here $\|\mathbf{A}\|_F$ denotes the Frobenius norm of the matrix \mathbf{A} viewed as a vector. Unfortunately, Powell's update does not preserve positive definiteness. (Hints: Note that $(\mathbf{N} - \mathbf{M}) \mathbf{d} = (\mathbf{N}_{\text{opt}} - \mathbf{M}) \mathbf{d}$ for every such \mathbf{N} and

$\|(\mathbf{N} - \mathbf{M})\mathbf{v}\| \geq \|(\mathbf{N}_{\text{opt}} - \mathbf{M})\mathbf{v}\|$ for every \mathbf{v} with $\mathbf{v}^* \mathbf{d} = 0$. Apply the identities

$$\|\mathbf{A}\|_F^2 = \|\mathbf{A}\mathbf{O}\|_F^2 = \sum_{i=1}^n \|\mathbf{A}\mathbf{o}_i\|_F^2$$

for an orthogonal matrix \mathbf{O} with columns $\mathbf{o}_1, \dots, \mathbf{o}_n$.)

8. Consider the quadratic function

$$Q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \begin{pmatrix} 2 & 1 \\ 1 & 1 \end{pmatrix} \mathbf{x} + (1, 1)\mathbf{x}$$

defined on \mathbb{R}^2 . Compute by hand the iterates of the conjugate gradient and BFGS algorithms starting from $\mathbf{x}_1 = \mathbf{0}$. For the BFGS algorithm take $\mathbf{H}_1 = \mathbf{I}$ and use an exact line search. You should find that the two sequences of iterates coincide. This phenomenon holds more generally for any strictly convex quadratic function in the BFGS algorithm given $\mathbf{H}_1 = \mathbf{I}$ [200].

9. Write a program to implement the conjugate gradient algorithm. Apply it to the function

$$f(\mathbf{x}) = \frac{1}{4}x_1^4 + \frac{1}{2}x_2^2 - x_1x_2 + x_1 - x_2$$

with two local minima. Demonstrate that your program will converge to either minimum depending on its starting value.

10. Prove Woodbury's generalization

$$(\mathbf{A} + \mathbf{U}\mathbf{B}\mathbf{V}^*)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{U}(\mathbf{B}^{-1} + \mathbf{V}^*\mathbf{A}^{-1}\mathbf{U})^{-1}\mathbf{V}^*\mathbf{A}^{-1}$$

of the Sherman-Morrison matrix inversion formula for compatible matrices \mathbf{A} , \mathbf{B} , \mathbf{U} , and \mathbf{V} . Apply this formula to the BFGS rank-two update.

11. A quasi-Newton minimization of the strictly convex quadratic function (11.2) generates a sequence of points $\mathbf{x}_1, \dots, \mathbf{x}_{n+1}$ with \mathbf{A} -conjugate differences $\mathbf{d}_i = \mathbf{x}_{i+1} - \mathbf{x}_i$. At the final iteration we have argued that the approximate Hessian satisfies $\mathbf{H}_{n+1}\mathbf{d}_i = \mathbf{g}_i$ for $1 \leq i \leq n$ and $\mathbf{g}_i = \nabla f(\mathbf{x}_{i+1}) - \nabla f(\mathbf{x}_i)$. Show that this implies $\mathbf{H}_{n+1} = \mathbf{A}$.
12. In the rank-one update, suppose we want both \mathbf{H}_{i+1} to remain positive definite and $\det \mathbf{H}_{i+1}$ to exceed some constant $\epsilon > 0$. Explain how these criteria can be simultaneously met by replacing $c_i < 0$ by

$$\max \left\{ c_i, \left(\frac{\epsilon}{\det \mathbf{H}_i} - 1 \right) \frac{1}{\mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i} \right\}$$

in updating \mathbf{H}_i . (Hint: In verifying this sufficient condition, you may want to use the one-dimensional version of the identity

$$\det(\mathbf{A}) \det(\mathbf{B}^{-1} - \mathbf{U}^* \mathbf{A}^{-1} \mathbf{U}) = \det(\mathbf{A} - \mathbf{U} \mathbf{B} \mathbf{U}^*) \det(\mathbf{B}^{-1})$$

for compatible matrices \mathbf{A} , \mathbf{B} , and \mathbf{U} .)

13. Let \mathbf{H} be a positive definite matrix. Prove [33] that

$$\operatorname{tr}(\mathbf{H}) - \ln \det(\mathbf{H}) \geq \ln[\operatorname{cond}_2(\mathbf{H})]. \quad (11.29)$$

The condition number $\operatorname{cond}_2(\mathbf{H})$ of \mathbf{H} equals $\|\mathbf{H}\| \cdot \|\mathbf{H}^{-1}\|$, that is the ratio of the largest to smallest eigenvalue of \mathbf{H} . (Hint: Express $\operatorname{tr}(\mathbf{H}) - \ln \det(\mathbf{H})$ in terms of the eigenvalues of \mathbf{H} . Then use the inequalities $\lambda - \ln \lambda \geq 1$ and $\lambda > 2 \ln \lambda$ for all $\lambda > 0$.)

14. In Davidon's symmetric rank-one update (11.15), it is possible to control the condition number of \mathbf{H}_{i+1} by shrinking the constant c_i . Suppose a moderately sized number δ is chosen. Due to inequality (11.29), one can avoid ill-conditioning in the matrices \mathbf{H}_i by imposing the constraint $\operatorname{tr}(\mathbf{H}_i) - \ln \det(\mathbf{H}_i) \leq \delta$. To see how this fits into the updating scheme (11.15), verify that

$$\begin{aligned} \ln \det(\mathbf{H}_{i+1}) &= \ln \det(\mathbf{H}_i) + \ln(1 + c_i \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i) \\ \operatorname{tr}(\mathbf{H}_{i+1}) &= \operatorname{tr}(\mathbf{H}_i) + c_i \|\mathbf{v}_i\|^2. \end{aligned}$$

Employing these results, deduce that $\operatorname{tr}(\mathbf{H}_{i+1}) - \ln \det(\mathbf{H}_{i+1}) \leq \delta$ provided c_i satisfies

$$c_i \|\mathbf{v}_i\|^2 - \ln(1 + c_i \mathbf{v}_i^* \mathbf{H}_i^{-1} \mathbf{v}_i) \leq \delta - \operatorname{tr}(\mathbf{H}_i) + \ln \det(\mathbf{H}_i).$$

15. Suppose the $n \times n$ symmetric matrix \mathbf{A} has eigenvalues

$$\lambda_1 < \lambda_2 \leq \cdots \leq \lambda_{n-1} < \lambda_n.$$

The iterative scheme $\mathbf{x}_{i+1} = (\mathbf{A} - \eta_i \mathbf{I}) \mathbf{x}_i$ can be used to approximate either λ_1 or λ_n . Consider the criterion

$$\sigma_i = \frac{\mathbf{x}_{i+1}^* \mathbf{A} \mathbf{x}_{i+1}}{\mathbf{x}_{i+1}^* \mathbf{x}_{i+1}}.$$

Choosing η_i to maximize σ_i causes $\lim_{i \rightarrow \infty} \sigma_i = \lambda_n$, while choosing η_i to minimize σ_i causes $\lim_{i \rightarrow \infty} \sigma_i = \lambda_1$. Show that the extrema of σ_i as a function of η are given by the roots of the quadratic equation

$$0 = \det \begin{pmatrix} 1 & \eta & \eta^2 \\ \tau_0 & \tau_1 & \tau_2 \\ \tau_1 & \tau_2 & \tau_3 \end{pmatrix},$$

where $\tau_k = \mathbf{x}_i^* \mathbf{A}^k \mathbf{x}_i$. Apply this algorithm to find the largest and smallest eigenvalue of the matrix

$$\mathbf{A} = \begin{pmatrix} 10 & 7 & 8 & 7 \\ 7 & 5 & 6 & 5 \\ 8 & 6 & 10 & 9 \\ 7 & 5 & 9 & 10 \end{pmatrix}.$$

You should find $\lambda_1 = 0.01015$ and $\lambda_4 = 30.2887$ [48].

16. Calculate the second derivative $\phi''(\mu)$ of the function defined in equation (11.27). Prove that $\phi''(\mu) \geq 0$.
17. Show that the solution $\mathbf{y}(\mu)$ of the equation $(\mathbf{H} + \mu\mathbf{I})\mathbf{y}(\mu) = \mathbf{e}$ satisfies $\lim_{\mu \rightarrow \infty} \|\mathbf{y}(\mu)\| = 0$. How does this justify the conclusion that the function (11.27) has a zero on $(0, \infty)$ when $\phi(0) > 0$?

12

Analysis of Convergence

12.1 Introduction

Proving convergence of the various optimization algorithms is a delicate exercise. In general, it is helpful to consider local and global convergence patterns separately. The local convergence rate of an algorithm provides a useful benchmark for comparing it to other algorithms. On this basis, Newton's method wins hands down. However, the tradeoffs are subtle. Besides the sheer number of iterations until convergence, the computational complexity and numerical stability of an algorithm are critically important. The MM algorithm is often the epitome of numerical stability and computational simplicity. Scoring lies somewhere between Newton's method and the MM algorithm. It tends to converge more quickly than the MM algorithm and to behave more stably than Newton's method. Quasi-Newton methods also occupy this intermediate zone. Because the issues are complex, all of these algorithms survive and prosper in certain computational niches.

The following short overview of convergence manages to cover only highlights. For the sake of simplicity, only unconstrained problems are treated. Quasi-Newton methods are also ignored. The efforts of a generation of numerical analysts in understanding quasi-Newton methods defy easy summary or digestion. Interested readers can consult one of the helpful references [66, 105, 183, 204]. We emphasize MM and gradient algorithms, partially because a fairly coherent theory for them can be reviewed in a few pages.

12.2 Local Convergence

Local convergence of many optimization algorithms hinges on the following result [207].

Proposition 12.2.1 (Ostrowski) *Let the differentiable map $h(\mathbf{x})$ from an open set $U \subset \mathbb{R}^n$ into \mathbb{R}^n have fixed point \mathbf{y} . If $\|dh(\mathbf{y})\|_{\dagger} < 1$ for some induced matrix norm, and if \mathbf{x}_0 is sufficiently close to \mathbf{y} , then the iterates $\mathbf{x}_{m+1} = h(\mathbf{x}_m)$ converge to \mathbf{y} .*

Proof: Let $h(\mathbf{x})$ have slope function $s(\mathbf{x}, \mathbf{y})$ near \mathbf{y} . For any constant r satisfying $\|dh(\mathbf{y})\|_{\dagger} < r < 1$, we have $\|s(\mathbf{x}, \mathbf{y})\|_{\dagger} \leq r$ for \mathbf{x} sufficiently close to \mathbf{y} . It therefore follows from the identities

$$\mathbf{x}_{m+1} - \mathbf{y} = h(\mathbf{x}_m) - h(\mathbf{y}) = s(\mathbf{x}_m, \mathbf{y})(\mathbf{x}_m - \mathbf{y})$$

that a proper choice of \mathbf{x}_0 yields

$$\|\mathbf{x}_{m+1} - \mathbf{y}\|_{\dagger} \leq \|s(\mathbf{x}_m, \mathbf{y})\|_{\dagger} \|\mathbf{x}_m - \mathbf{y}\|_{\dagger} \leq r \|\mathbf{x}_m - \mathbf{y}\|_{\dagger}.$$

In other words, the distance from \mathbf{x}_m to \mathbf{y} contracts by a factor of at least r at every iteration. This proves convergence. ■

Two comments are worth making about Proposition 12.2.1. First, the appearance of a general vector norm and its induced matrix norm obscures the fact that the condition $\rho[dh(\mathbf{y})] < 1$ on the spectral radius of $dh(\mathbf{y})$ is the operative criterion. One can prove that any induced matrix norm exceeds the spectral radius and that some induced matrix norm comes within ϵ of it for any small $\epsilon > 0$ [166]. Later in this section, we will generate a tight matrix norm by taking an $n \times n$ invertible matrix \mathbf{T} and forming $\|\mathbf{u}\|_{\mathbf{T}} = \|\mathbf{T}\mathbf{u}\|$. It is easy to check that this defines a legitimate vector norm and that the induced matrix norm $\|\mathbf{M}\|_{\mathbf{T}}$ on $n \times n$ matrices \mathbf{M} satisfies

$$\|\mathbf{M}\|_{\mathbf{T}} = \sup_{\mathbf{u} \neq \mathbf{0}} \frac{\|\mathbf{T}\mathbf{M}\mathbf{u}\|}{\|\mathbf{T}\mathbf{u}\|} = \sup_{\mathbf{v} \neq \mathbf{0}} \frac{\|\mathbf{T}\mathbf{M}\mathbf{T}^{-1}\mathbf{v}\|}{\|\mathbf{v}\|}.$$

In other words, $\|\mathbf{M}\|_{\mathbf{T}} = \|\mathbf{T}\mathbf{M}\mathbf{T}^{-1}\|$, and we are back in the familiar terrain covered by the spectral norm.

Our second comment involves two definitions. A sequence \mathbf{x}_m is said to converge linearly to a point \mathbf{y} at rate $r < 1$ provided

$$\lim_{m \rightarrow \infty} \frac{\|\mathbf{x}_{m+1} - \mathbf{y}\|}{\|\mathbf{x}_m - \mathbf{y}\|} = r.$$

The sequence converges quadratically if the limit

$$\lim_{m \rightarrow \infty} \frac{\|\mathbf{x}_{m+1} - \mathbf{y}\|}{\|\mathbf{x}_m - \mathbf{y}\|^2} = c$$

exists. Ostrowski's result guarantees at least linear convergence; Newton's method improves linear convergence to quadratic convergence.

Our intention is to apply Ostrowski's result to iteration maps of the type

$$h(\mathbf{x}) = \mathbf{x} - A(\mathbf{x})^{-1}b(\mathbf{x}). \quad (12.1)$$

A point \mathbf{y} is fixed by the map $h(\mathbf{x})$ if and only if $b(\mathbf{y}) = \mathbf{0}$. In optimization problems, $b(\mathbf{x}) = \nabla f(\mathbf{x})$ for some real-valued function $f(\mathbf{x})$ defined on \mathbb{R}^n . Thus, fixed points correspond to stationary points. The matrix $A(\mathbf{x})$ is typically $d^2f(\mathbf{x})$ or a positive definite or negative definite approximation to it. For instance, in statistical applications, $-A(\mathbf{x})$ could be either the observed or expected information. In the MM gradient algorithm, $A(\mathbf{x})$ is the second differential $d^2g(\mathbf{x} \mid \mathbf{x}_m)$ of the surrogate function.

Our first order of business is to compute the differential $dh(\mathbf{y})$ and an associated slope function $s_h(\mathbf{x}, \mathbf{y})$ at a fixed point \mathbf{y} of $h(\mathbf{x})$ in terms of the slope function $s_b(\mathbf{x}, \mathbf{y})$ of $b(\mathbf{x})$. Because $b(\mathbf{y}) = \mathbf{0}$ at a fixed point, the calculation

$$\begin{aligned} h(\mathbf{x}) - h(\mathbf{y}) &= \mathbf{x} - \mathbf{y} - A(\mathbf{x})^{-1}[b(\mathbf{x}) - b(\mathbf{y})] \\ &= [\mathbf{I} - A(\mathbf{x})^{-1}s_b(\mathbf{x}, \mathbf{y})](\mathbf{x} - \mathbf{y}) \end{aligned} \quad (12.2)$$

identifies the slope function

$$s_h(\mathbf{x}, \mathbf{y}) = \mathbf{I} - A(\mathbf{x})^{-1}s_b(\mathbf{x}, \mathbf{y})$$

and corresponding differential

$$dh(\mathbf{y}) = \mathbf{I} - A(\mathbf{y})^{-1}db(\mathbf{y}). \quad (12.3)$$

In Newton's method, $A(\mathbf{y}) = db(\mathbf{y})$ and

$$\mathbf{I} - A(\mathbf{y})^{-1}db(\mathbf{y}) = \mathbf{I} - db(\mathbf{y})^{-1}db(\mathbf{y}) = \mathbf{0}.$$

Proposition 12.2.1 therefore implies that the Newton iterates are locally attracted to a fixed point \mathbf{y} . Of course, this conclusion is predicated on the suppositions that $db(\mathbf{x})$ tends to $db(\mathbf{y})$ as \mathbf{x} tends to \mathbf{y} and that $db(\mathbf{y})$ is invertible. To demonstrate quadratic rather than linear convergence, we now assume that $b(\mathbf{x})$ is differentiable and that its differential $db(\mathbf{x})$ is Lipschitz in a neighborhood of \mathbf{y} with Lipschitz constant λ . Given these assumptions and the identities $h(\mathbf{y}) = \mathbf{y}$ and $A(\mathbf{x}) = db(\mathbf{x})$, equation (12.2) implies

$$\begin{aligned} h(\mathbf{x}) - \mathbf{y} &= -db(\mathbf{x})^{-1}[b(\mathbf{x}) - b(\mathbf{y}) - db(\mathbf{y})(\mathbf{x} - \mathbf{y})] \\ &\quad + db(\mathbf{x})^{-1}[db(\mathbf{x}) - db(\mathbf{y})](\mathbf{x} - \mathbf{y}) \end{aligned}$$

Since Problem 31 of Chap. 4 supplies the bound

$$\|b(\mathbf{x}) - b(\mathbf{y}) - db(\mathbf{y})(\mathbf{x} - \mathbf{y})\| \leq \frac{\lambda}{2}\|\mathbf{x} - \mathbf{y}\|^2,$$

it follows that

$$\|h(\mathbf{x}) - \mathbf{y}\| \leq \left(\frac{\lambda}{2} + \lambda\right) \|db(\mathbf{x})^{-1}\| \cdot \|\mathbf{x} - \mathbf{y}\|^2.$$

The next proposition summarizes our results.

Proposition 12.2.2 *Let \mathbf{y} be a $\mathbf{0}$ of the continuously differentiable function $b(\mathbf{x})$ from an open set $U \subset \mathbb{R}^n$ into \mathbb{R}^n . If $db(\mathbf{y})^{-1}$ is invertible and $db(\mathbf{x})$ is Lipschitz in U with constant λ , then Newton's method converges to \mathbf{y} at a quadratic rate or better whenever \mathbf{x}_0 is sufficiently close to \mathbf{y} .*

Proof: The preceding remarks make it clear that

$$\limsup_{m \rightarrow \infty} \frac{\|\mathbf{x}_{m+1} - \mathbf{y}\|}{\|\mathbf{x}_m - \mathbf{y}\|^2} \leq \frac{3\lambda}{2} \|db(\mathbf{y})^{-1}\|,$$

and this suffices for quadratic convergence or better. ■

We now turn to the MM gradient algorithm. Suppose we are minimizing $f(\mathbf{x})$ via the surrogate function $g(\mathbf{x} | \mathbf{x}_m)$. If \mathbf{y} is a local minimum of $f(\mathbf{x})$, it is reasonable to assume that the matrices $\mathbf{C} = d^2f(\mathbf{y})$ and $\mathbf{D} = d^2g(\mathbf{y} | \mathbf{y})$ are positive definite. Because $g(\mathbf{x} | \mathbf{y}) - f(\mathbf{x})$ attains its minimum at $\mathbf{x} = \mathbf{y}$, the matrix difference $\mathbf{D} - \mathbf{C}$ is certainly positive semidefinite. The MM gradient algorithm iterates take $A(\mathbf{x}) = d^2g(\mathbf{x} | \mathbf{x})$. In view of formula (12.3), the iteration map $h(\mathbf{x})$ has differential $\mathbf{I} - \mathbf{D}^{-1}\mathbf{C}$ at \mathbf{y} . If we let \mathbf{T} be the symmetric square root $\mathbf{D}^{1/2}$ of \mathbf{D} , then

$$\begin{aligned} \mathbf{I} - \mathbf{D}^{-1}\mathbf{C} &= \mathbf{D}^{-1}(\mathbf{D} - \mathbf{C}) \\ &= \mathbf{T}^{-1}\mathbf{T}^{-1}(\mathbf{D} - \mathbf{C})\mathbf{T}^{-1}\mathbf{T}. \end{aligned}$$

Hence, $\mathbf{I} - \mathbf{D}^{-1}\mathbf{C}$ is similar to $\mathbf{T}^{-1}(\mathbf{D} - \mathbf{C})\mathbf{T}^{-1}$.

To establish local attraction of the MM gradient algorithm to \mathbf{y} , we need to choose an appropriate matrix norm. The choice $\|\mathbf{M}\|_{\mathbf{T}} = \|\mathbf{T}\mathbf{M}\mathbf{T}^{-1}\|$ serves well because Example 1.4.3 and Proposition 2.2.1 imply that

$$\begin{aligned} \|\mathbf{I} - \mathbf{D}^{-1}\mathbf{C}\|_{\mathbf{T}} &= \|\mathbf{T}\mathbf{T}^{-1}\mathbf{T}^{-1}(\mathbf{D} - \mathbf{C})\mathbf{T}^{-1}\mathbf{T}\mathbf{T}^{-1}\| \\ &= \|\mathbf{T}^{-1}(\mathbf{D} - \mathbf{C})\mathbf{T}^{-1}\| \\ &= \sup_{\mathbf{u} \neq \mathbf{0}} \frac{\mathbf{u}^*\mathbf{T}^{-1}(\mathbf{D} - \mathbf{C})\mathbf{T}^{-1}\mathbf{u}}{\mathbf{u}^*\mathbf{u}} \\ &= \sup_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^*(\mathbf{D} - \mathbf{C})\mathbf{v}}{\mathbf{v}^*\mathbf{T}^*\mathbf{T}\mathbf{v}} \\ &= \sup_{\mathbf{v} \neq \mathbf{0}} \frac{\mathbf{v}^*(\mathbf{D} - \mathbf{C})\mathbf{v}}{\mathbf{v}^*\mathbf{D}\mathbf{v}} \\ &= 1 - \inf_{\|\mathbf{v}\|=1} \frac{\mathbf{v}^*\mathbf{C}\mathbf{v}}{\mathbf{v}^*\mathbf{D}\mathbf{v}}. \end{aligned}$$

The symmetry and positive semidefiniteness of $\mathbf{D} - \mathbf{C}$ come into play in the third equality in this string of equalities. By virtue of the positive definiteness of \mathbf{C} and \mathbf{D} , the continuous ratio $\mathbf{v}^* \mathbf{C} \mathbf{v} / \mathbf{v}^* \mathbf{D} \mathbf{v}$ is bounded below by a positive constant on the compact sphere $\{\mathbf{v} : \|\mathbf{v}\| = 1\}$. It follows that $\|\mathbf{I} - \mathbf{D}^{-1} \mathbf{C}\|_{\mathcal{T}} < 1$, and Ostrowski's result applies. Hence, the iterates \mathbf{x}_m are locally attracted to \mathbf{y} .

Calculation of the differential $dh(\mathbf{y})$ of an MM iteration map $h(\mathbf{x})$ is equally interesting. This map satisfies the equation

$$\nabla g[h(\mathbf{x}) \mid \mathbf{x}] = \mathbf{0}$$

Assuming that the matrix $d^2g(\mathbf{y} \mid \mathbf{y})$ is invertible, the implicit function theorem, Proposition 4.6.2, shows that $h(\mathbf{x})$ is continuously differentiable with differential

$$dh(\mathbf{x}) = -d^2g[h(\mathbf{x}) \mid \mathbf{x}]^{-1} d^{11}g[h(\mathbf{x}) \mid \mathbf{x}]. \quad (12.4)$$

Here $d^{11}g(\mathbf{u} \mid \mathbf{v})$ denotes the differential of $dg(\mathbf{u} \mid \mathbf{v})$ with respect to \mathbf{v} . At the fixed point \mathbf{y} of $h(\mathbf{x})$, equation (12.4) becomes

$$dh(\mathbf{y}) = -d^2g(\mathbf{y} \mid \mathbf{y})^{-1} d^{11}g(\mathbf{y} \mid \mathbf{y}). \quad (12.5)$$

Further simplification can be achieved by taking the differential of

$$\nabla f(\mathbf{x}) - \nabla g(\mathbf{x} \mid \mathbf{x}) = \mathbf{0}$$

and setting $\mathbf{x} = \mathbf{y}$. These actions give

$$d^2f(\mathbf{y}) - d^2g(\mathbf{y} \mid \mathbf{y}) - d^{11}g(\mathbf{y} \mid \mathbf{y}) = \mathbf{0}.$$

This last equation can be solved for $d^{11}g(\mathbf{y} \mid \mathbf{y})$, and the result substituted in equation (12.5). It follows that

$$\begin{aligned} dh(\mathbf{y}) &= -d^2g(\mathbf{y} \mid \mathbf{y})^{-1} [d^2f(\mathbf{y}) - d^2g(\mathbf{y} \mid \mathbf{y})] \\ &= \mathbf{I} - d^2g(\mathbf{y} \mid \mathbf{y})^{-1} d^2f(\mathbf{y}), \end{aligned} \quad (12.6)$$

which is precisely the differential computed for the MM gradient algorithm. Hence, the MM and MM gradient algorithms display exactly the same behavior in converging to a stationary point of $f(\mathbf{x})$.

These apparently esoteric details have considerable practical value. In the setting of the EM algorithm, we replace \mathbf{x} by $\boldsymbol{\theta}$, $f(\mathbf{x})$ by the observed data loglikelihood $L(\boldsymbol{\theta})$, and $g(\mathbf{x} \mid \mathbf{x}_m)$ by the minorizing function $Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$. The rate of convergence of an EM algorithm is determined by how well the observed data information $-d^2L(\boldsymbol{\theta})$ is approximated by the complete data information matrix $-d^2Q(\boldsymbol{\theta} \mid \boldsymbol{\theta})$ at the optimal point. The difference between these two matrices is termed the missing data information [191, 206]. If the amount of missing data is high, then the algorithm will converge

slowly. The art in devising an EM or MM algorithm lies in choosing a tractable surrogate that matches the objective function as closely possible.

Local convergence of the scoring algorithm is not guaranteed by Proposition 12.2.1 because nothing prevents an eigenvalue of

$$dh(\mathbf{y}) = \mathbf{I} + J(\mathbf{y})^{-1}d^2L(\mathbf{y})$$

from falling below -1 . Here $L(\mathbf{x})$ is the loglikelihood, $J(\mathbf{x})$ is the expected information, and $h(\mathbf{x})$ is the scoring iteration map. Scoring with a fixed partial step,

$$\mathbf{x}_{m+1} = \mathbf{x}_m + tJ(\mathbf{x}_m)^{-1}\nabla L(\mathbf{x}_m),$$

will converge locally for $t > 0$ sufficiently small. In practice, no adjustment is usually necessary. For reasonably large sample sizes, the expected information matrix $J(\mathbf{y})$ approximates the observed information matrix $-d^2L(\mathbf{y})$ well, and the spectral radius of $dh(\mathbf{y})$ is nearly 0.

Finally, let us consider local convergence of block relaxation. The argument $\mathbf{x} = (\mathbf{x}_{[1]}, \mathbf{x}_{[2]}, \dots, \mathbf{x}_{[b]})$ of the objective function $f(\mathbf{x})$ now splits into disjoint blocks, and $f(\mathbf{x})$ is minimized along each block of components $\mathbf{x}_{[i]}$ in turn. Let $M_i(\mathbf{x})$ denote the update to block i . To compute the differential of the full update $M(\mathbf{x})$ at a local optimum \mathbf{y} , we need compact notation. Let $\partial_i f(\mathbf{x})$ denote the partial differential of $f(\mathbf{x})$ with respect to block i ; the transpose of $\partial_i f(\mathbf{x})$ is the partial gradient $\nabla_i f(\mathbf{x})$. The updates satisfy the partial gradient equations

$$\mathbf{0} = \nabla_i f[M_1(\mathbf{x}), \dots, M_i(\mathbf{x}), \mathbf{x}_{[i+1]}, \dots, \mathbf{x}_{[b]}]. \quad (12.7)$$

Now let $\partial_j \nabla_i f(\mathbf{x})$ denote the partial differential of the partial gradient $\nabla_i f(\mathbf{x})$ with respect to block j . Taking the partial differential of equation (12.7) with respect to block j , applying the chain rule, and substituting the optimal point $\mathbf{y} = M(\mathbf{y})$ for \mathbf{x} yield

$$\begin{aligned} \mathbf{0} &= \sum_{k=1}^i \partial_k \nabla_i f(\mathbf{y}) \partial_j M_k(\mathbf{y}), \quad j \leq i \\ \mathbf{0} &= \sum_{k=1}^i \partial_k \nabla_i f(\mathbf{y}) \partial_j M_k(\mathbf{y}) + \partial_j \nabla_i f(\mathbf{y}), \quad j > i. \end{aligned} \quad (12.8)$$

It is helpful to express these equations in block matrix form.

For example in the case of $b = 3$ blocks, the linear system of equations (12.8) can be represented as $\mathbf{L}dM(\mathbf{y}) = \mathbf{D} - \mathbf{U}$, where $\mathbf{U} = \mathbf{L}^*$ and

$$dM(\mathbf{y}) = \begin{pmatrix} \partial_1 M_1(\mathbf{y}) & \partial_2 M_1(\mathbf{y}) & \partial_3 M_1(\mathbf{y}) \\ \partial_1 M_2(\mathbf{y}) & \partial_2 M_2(\mathbf{y}) & \partial_3 M_2(\mathbf{y}) \\ \partial_1 M_3(\mathbf{y}) & \partial_2 M_3(\mathbf{y}) & \partial_3 M_3(\mathbf{y}) \end{pmatrix}$$

$$\mathbf{L} = \begin{pmatrix} \partial_1 \nabla_1 f(\mathbf{y}) & \mathbf{0} & \mathbf{0} \\ \partial_1 \nabla_2 f(\mathbf{y}) & \partial_2 \nabla_2 f(\mathbf{y}) & \mathbf{0} \\ \partial_1 \nabla_3 f(\mathbf{y}) & \partial_2 \nabla_3 f(\mathbf{y}) & \partial_3 \nabla_3 f(\mathbf{y}) \end{pmatrix}$$

$$\mathbf{D} = \begin{pmatrix} \partial_1 \nabla_1 f(\mathbf{y}) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \partial_2 \nabla_2 f(\mathbf{y}) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \partial_3 \nabla_3 f(\mathbf{y}) \end{pmatrix}.$$

The identity $[\partial_j \nabla_i f(\mathbf{y})]^* = \partial_i \nabla_j f(\mathbf{y})$ between two nontrivial blocks of \mathbf{U} and \mathbf{L} is a consequence of the equality of mixed partials. The matrix equation $\mathbf{L} dM(\mathbf{y}) = \mathbf{D} - \mathbf{U}$ can be explicitly solved in the form

$$dM(\mathbf{y}) = \mathbf{L}^{-1}(\mathbf{D} - \mathbf{U}).$$

Here \mathbf{L} is invertible provided its diagonal blocks $\partial_i \nabla_i f(\mathbf{y})$ are invertible. At an optimal point \mathbf{y} , the partial Hessian matrix $\partial_i \nabla_i f(\mathbf{y})$ is always positive semidefinite and usually positive definite as well.

Local convergence of block relaxation hinges on whether the spectral radius ρ of the matrix $\mathbf{L}^{-1}(\mathbf{U} - \mathbf{D})$ satisfies $\rho < 1$. Suppose that λ is an eigenvalue of $\mathbf{L}^{-1}(\mathbf{D} - \mathbf{U})$ with eigenvector \mathbf{v} . These can be complex. The equality $\mathbf{L}^{-1}(\mathbf{D} - \mathbf{U})\mathbf{v} = \lambda\mathbf{v}$ implies $(1 - \lambda)\mathbf{L}\mathbf{v} = (\mathbf{L} + \mathbf{U} - \mathbf{D})\mathbf{v}$. Premultiplying this by the conjugate transpose \mathbf{v}^* gives

$$\frac{1}{1 - \lambda} = \frac{\mathbf{v}^* \mathbf{L} \mathbf{v}}{\mathbf{v}^* (\mathbf{L} + \mathbf{U} - \mathbf{D}) \mathbf{v}}.$$

Hence, the real part of $1/(1 - \lambda)$ satisfies

$$\begin{aligned} \operatorname{Re}\left(\frac{1}{1 - \lambda}\right) &= \frac{\mathbf{v}^* (\mathbf{L} + \mathbf{U}) \mathbf{v}}{2\mathbf{v}^* (\mathbf{L} + \mathbf{U} - \mathbf{D}) \mathbf{v}} \\ &= \frac{1}{2} \left[1 + \frac{\mathbf{v}^* \mathbf{D} \mathbf{v}}{\mathbf{v}^* d^2 f(\mathbf{y}) \mathbf{v}} \right] \\ &> \frac{1}{2} \end{aligned}$$

for $d^2 f(\mathbf{y})$ positive definite. If $\lambda = \alpha + \beta\sqrt{-1}$, then the last inequality entails

$$\frac{1 - \alpha}{(1 - \alpha)^2 + \beta^2} > \frac{1}{2},$$

which is equivalent to $|\lambda|^2 = \alpha^2 + \beta^2 < 1$. Hence, the spectral radius $\rho < 1$.

12.3 Coercive Functions

The concept of coerciveness is critical in establishing the existence of minimum points. A function $f(\mathbf{x})$ on \mathbb{R}^n is said to be coercive if

$$\lim_{\|\mathbf{x}\| \rightarrow \infty} f(\mathbf{x}) = \infty.$$

If $f(\mathbf{x})$ is both lower semicontinuous and coercive, then all sublevel sets $\{\mathbf{x} : f(\mathbf{x}) \leq c\}$ are compact, and the minimum value of $f(\mathbf{x})$ is attained. This improvement of Weierstrass' theorem (Proposition 2.5.4) plays a key role in optimization theory.

Two strategies stand out in proving coerciveness. One revolves around comparing one function to another. For example, suppose $f(x) = x^2 + \sin x$ and $g(x) = x^2 - 1$. Then $g(x)$ is clearly coercive and $f(x) \geq g(x)$. Hence, $f(x)$ is also coercive. As explained in the next proposition, the second strategy is restricted to convex functions. In stating the proposition, we allow $f(\mathbf{x})$ to have the value ∞ .

Proposition 12.3.1 *Suppose $f(\mathbf{x})$ is a convex lower semicontinuous function on \mathbb{R}^n . Choose any point \mathbf{y} with $f(\mathbf{y}) < \infty$. Then $f(\mathbf{x})$ is coercive if and only if $f(\mathbf{x})$ is coercive along all nontrivial rays $\{\mathbf{x} \in \mathbb{R}^n : \mathbf{x} = \mathbf{y} + t\mathbf{v}, t \geq 0\}$ emanating from \mathbf{y} .*

Proof: The stated condition is obviously necessary. To prove that it is sufficient, suppose it holds, but $f(\mathbf{x})$ is not coercive. It suffices to take $\mathbf{y} = \mathbf{0}$ because one can always consider the translated function $g(\mathbf{x}) = f(\mathbf{x} - \mathbf{y})$, which retains the properties of convexity and lower semicontinuity. Let \mathbf{x}_n be a sequence such that $\lim_{n \rightarrow \infty} \|\mathbf{x}_n\| = \infty$ and $\limsup_{n \rightarrow \infty} f(\mathbf{x}_n) < \infty$. By passing to a subsequence if necessary, we can assume that the unit vectors $\mathbf{v}_n = \|\mathbf{x}_n\|^{-1} \mathbf{x}_n$ converge to a unit vector \mathbf{v} . For $t > 0$ and n large enough, convexity implies

$$f(t\mathbf{v}_n) \leq \frac{t}{\|\mathbf{x}_n\|} f(\mathbf{x}_n) + \left(1 - \frac{t}{\|\mathbf{x}_n\|}\right) f(\mathbf{0}).$$

It follows that $\liminf_{n \rightarrow \infty} f(t\mathbf{v}_n) \leq f(\mathbf{0})$. On the other hand, lower semicontinuity entails

$$f(t\mathbf{v}) \leq \liminf_{n \rightarrow \infty} f(t\mathbf{v}_n) \leq f(\mathbf{0}).$$

Hence, $f(\mathbf{x})$ does not tend to ∞ along the ray $t\mathbf{v}$, contradicting our assumption. ■

For example, consider the quadratic $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$. If \mathbf{A} is positive definite and $\mathbf{v} \neq \mathbf{0}$, then $f(t\mathbf{v})$ is a quadratic in t with positive leading coefficient. Hence, $f(\mathbf{x})$ is coercive along the ray $t\mathbf{v}$. It follows that $f(\mathbf{x})$ is coercive. When \mathbf{A} is positive semidefinite and $\mathbf{A}\mathbf{v} = \mathbf{0}$, $f(t\mathbf{v})$ is linear in t . If the leading coefficient $\mathbf{b}^* \mathbf{v}$ is positive, then $f(\mathbf{x})$ is coercive along the ray $t\mathbf{v}$. However, $f(\mathbf{x})$ then fails to be coercive along the ray $-\mathbf{v}$. Hence, $f(\mathbf{x})$ is not coercive. This fact does not prevent $f(\mathbf{x})$ from attaining its minimum. If \mathbf{x} satisfies $\mathbf{A}\mathbf{x} = -\mathbf{b}$, then the convex function $f(\mathbf{x})$ has a stationary point, which necessarily furnishes a minimum value.

Posynomials present another interesting test case. In the exponential parameterization $t_k = e^{x_k}$, one can represent a posynomial as

$$f(\mathbf{x}) = \sum_{i=1}^j c_i e^{\beta_i^* \mathbf{x}}.$$

At the preferred point $\mathbf{0}$, it is obvious that $f(t\mathbf{v})$ tends to ∞ as t tends to ∞ if and only if at least one β_i satisfies $\beta_i^* \mathbf{v} > 0$. In other words, $f(\mathbf{x})$ is coercive if and only if the polar cone

$$C = \{\mathbf{v} : \beta_i^* \mathbf{v} \leq 0 \text{ for all } i\}$$

consists of the trivial vector $\mathbf{0}$ alone. Section 14.3.7 treats polar cones in more detail.

In the original posynomial parameterization, t_k tends to 0 as x_k tends to $-\infty$. This suggests the need for a broader definition of coerciveness consistent with Weierstrass' theorem. Suppose the lower semicontinuous function $f(\mathbf{x})$ is defined on an open set U . To avoid colliding with the boundary of U , we assume that the set

$$C_{\mathbf{y}} = \{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{y})\}$$

is compact for every $\mathbf{y} \in U$. If this is the case, then $f(\mathbf{x})$ attains its minimum somewhere in U . The essence of the expanded definition of coerciveness is that $f(\mathbf{x})$ tends to ∞ as \mathbf{x} approaches the boundary of U or $\|\mathbf{x}\|$ approaches ∞ .

12.4 Global Convergence of the MM Algorithm

In this section and the next, we tackle global convergence. We begin with the MM algorithm and consider without loss of generality minimization of the objective function $f(\mathbf{x})$ via the majorizing surrogate $g(\mathbf{x} \mid \mathbf{x}_m)$. In studying global convergence, we must carefully specify the parameter domain U . Let us take U to be any open convex subset of \mathbb{R}^n . It is convenient to assume that $f(\mathbf{x})$ is coercive on U in the sense just specified and that whenever necessary $f(\mathbf{x})$ and $g(\mathbf{x} \mid \mathbf{x}_m)$ and their various first and second differentials are jointly continuous in \mathbf{x} and \mathbf{x}_m .

We also demand that the second differential $d^2g(\mathbf{x} \mid \mathbf{x}_m)$ be positive definite. This implies that $g(\mathbf{x} \mid \mathbf{x}_m)$ is strictly convex. Strict convexity in turn implies that the solution \mathbf{x}_{m+1} of the minimization step is unique. Existence of a solution fortunately is guaranteed by coerciveness. Indeed, the closed set

$$\{\mathbf{x} : g(\mathbf{x} \mid \mathbf{x}_m) \leq g(\mathbf{x}_m \mid \mathbf{x}_m) = f(\mathbf{x}_m)\}$$

is compact because it is contained within the compact set

$$\{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{x}_m)\}.$$

Finally, the implicit function theorem, Proposition 4.6.2, shows that the iteration map $\mathbf{x}_{m+1} = M(\mathbf{x}_m)$ is continuously differentiable in a neighborhood of every point \mathbf{x}_m . Local differentiability of $M(\mathbf{x})$ clearly extends to global differentiability.

Gradient versions of the algorithm (12.1) have the property that stationary points of the objective function and fixed points of the iteration map coincide. This property also applies to the MM algorithm. Here we recall the two identities $\nabla g(\mathbf{x}_{m+1} | \mathbf{x}_m) = \mathbf{0}$ and $\nabla g(x_m | x_m) = \nabla f(\mathbf{x}_m)$ and the strict convexity of $g(\mathbf{x} | \mathbf{x}_m)$. By the same token, stationary points and only stationary points give equality in the descent inequality $f[M(\mathbf{x})] \leq f(\mathbf{x})$.

The next technical proposition prepares the ground for a proof of global convergence. We remind the reader that a point \mathbf{y} is a cluster point of a sequence \mathbf{x}_m provided there is a subsequence \mathbf{x}_{m_k} that tends to \mathbf{y} . One can easily verify that any limit of a sequence of cluster points is also a cluster point and that a bounded sequence has a limit if and only if it has at most one cluster point. See Problem 21.

Proposition 12.4.1 *If a bounded sequence \mathbf{x}_m in \mathbb{R}^n satisfies*

$$\lim_{m \rightarrow \infty} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| = 0, \quad (12.9)$$

then its set T of cluster points is connected. If T is finite, then T reduces to a single point, and $\lim_{m \rightarrow \infty} \mathbf{x}_m = \mathbf{y}$ exists.

Proof: It is straightforward to prove that T is a compact set. If it is disconnected, then there is a continuous disconnecting function $\phi(\mathbf{x})$ having exactly the two values 0 and 1. The inverse images of the closed sets 0 and 1 under $\phi(\mathbf{x})$ can be represented as the intersections $T_0 = T \cap C_0$ and $T_1 = T \cap C_1$ of T with two closed sets C_0 and C_1 . Because T is compact, T_0 and T_1 are closed, nonempty, and disjoint. Furthermore, the distance

$$\text{dist}(T_0, T_1) = \inf_{\mathbf{u} \in T_0} \text{dist}(\mathbf{u}, T_1) = \inf_{\mathbf{u} \in T_0, \mathbf{v} \in T_1} \|\mathbf{u} - \mathbf{v}\|$$

separating T_0 and T_1 is positive. Indeed, the continuous function $\text{dist}(\mathbf{u}, T_1)$ attains its minimum at some point \mathbf{u} of the compact set T_0 , and the distance $\text{dist}(\mathbf{u}, T_1)$ separating that \mathbf{u} from T_1 must be positive because T_1 is closed.

Now consider the sequence \mathbf{x}_m in the statement of the proposition. For large enough m , we have $\|\mathbf{x}_{m+1} - \mathbf{x}_m\| < \text{dist}(T_0, T_1)/4$. As the sequence \mathbf{x}_m bounces back and forth between cluster points in T_0 and T_1 , it must enter the closed set $W = \{\mathbf{u} : \text{dist}(\mathbf{u}, T) \geq \text{dist}(T_0, T_1)/4\}$ infinitely often. But this means that W contains a cluster point of \mathbf{x}_m . Because W is disjoint from T_0 and T_1 , and these two sets are postulated to contain all of the cluster points of \mathbf{x}_m , this contradiction implies that T is connected.

Because a finite set with more than one point is necessarily disconnected, T can be a finite set only if it consists of a single point. Finally, a bounded sequence with only a single cluster point has that point as its limit. ■

With these facts in mind, we now state and prove a version of Liapunov's theorem for discrete dynamical systems [183].

Proposition 12.4.2 (Liapunov) *Let Γ be the set of cluster points generated by the MM sequence $\mathbf{x}_{m+1} = M(\mathbf{x}_m)$ starting from some initial \mathbf{x}_0 . Then Γ is contained in the set S of stationary points of $f(\mathbf{x})$.*

Proof: The sequence \mathbf{x}_m stays within the compact set

$$\{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}.$$

Consider a cluster point $\mathbf{z} = \lim_{k \rightarrow \infty} \mathbf{x}_{m_k}$. Since the sequence $f(\mathbf{x}_m)$ is monotonically decreasing and bounded below, $\lim_{m \rightarrow \infty} f(\mathbf{x}_m)$ exists. Hence, taking limits in the inequality $f[M(\mathbf{x}_{m_k})] \leq f(\mathbf{x}_{m_k})$ and invoking the continuity of $M(\mathbf{x})$ and $f(\mathbf{x})$ imply $f[M(\mathbf{z})] = f(\mathbf{z})$. Thus, \mathbf{z} is a fixed point of $M(\mathbf{x})$ and consequently also a stationary point of $f(\mathbf{x})$. ■

The next two propositions are adapted from the reference [195]. In the second of these, recall that a point \mathbf{x} in a set S is isolated if and only if there exists a radius $r > 0$ such that $S \cap B(\mathbf{x}, r) = \{\mathbf{x}\}$.

Proposition 12.4.3 *The set of cluster points Γ of $\mathbf{x}_{m+1} = M(\mathbf{x}_m)$ is compact and connected.*

Proof: Γ is a closed subset of the compact set $\{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ and is therefore itself compact. According to Proposition 12.4.1, Γ is connected provided $\lim_{m \rightarrow \infty} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| = 0$. If this sufficient condition fails, then the compactness of $\{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ makes it possible to extract a subsequence \mathbf{x}_{m_k} such that $\lim_{k \rightarrow \infty} \mathbf{x}_{m_k} = \mathbf{u}$ and $\lim_{k \rightarrow \infty} \mathbf{x}_{m_k+1} = \mathbf{v}$ both exist, but $\mathbf{v} \neq \mathbf{u}$. However, the continuity of $M(\mathbf{x})$ requires $\mathbf{v} = M(\mathbf{u})$ while the descent condition implies

$$f(\mathbf{v}) = f(\mathbf{u}) = \lim_{m \rightarrow \infty} f(\mathbf{x}_m).$$

The equality $f(\mathbf{v}) = f(\mathbf{u})$ forces the contradictory conclusion that \mathbf{u} is a fixed point of $M(\mathbf{x})$. Hence, the sufficient condition (12.9) for connectivity holds. ■

Proposition 12.4.4 *Suppose that all stationary points of $f(\mathbf{x})$ are isolated and that the stated differentiability, coerciveness, and convexity assumptions are true. Then any sequence of iterates $\mathbf{x}_{m+1} = M(\mathbf{x}_m)$ generated by the iteration map $M(\mathbf{x})$ of the MM algorithm possesses a limit, and that limit is a stationary point of $f(\mathbf{x})$. If $f(\mathbf{x})$ is strictly convex, then $\lim_{m \rightarrow \infty} \mathbf{x}_m$ is the minimum point.*

Proof: In the compact set $\{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ there can only be a finite number of stationary points. An infinite number of stationary points would admit a convergent sequence whose limit would not be isolated. Since the set of cluster points Γ is a connected subset of this finite set of stationary points, Γ reduces to a single point. ■

Two remarks on Proposition 12.4.4 are in order. First, except when strict convexity prevails for $f(\mathbf{x})$, the proposition offers no guarantee that the limit \mathbf{y} of the sequence \mathbf{x}_m furnishes a global minimum. Problem 11 contains a counterexample of Wu [275] exhibiting convergence to a saddle point in the EM algorithm. Fortunately, in practice, descent algorithms almost always converge to at least a local minimum of the objective function. Second, suppose that the set S of stationary points possesses a sequence $\mathbf{z}_m \in S$ converging to $\mathbf{z} \in S$ with $\mathbf{z}_m \neq \mathbf{z}$ for all m . Because the unit sphere in \mathbb{R}^n is compact, we can extract a subsequence such that

$$\lim_{k \rightarrow \infty} \frac{\mathbf{z}_{m_k} - \mathbf{z}}{\|\mathbf{z}_{m_k} - \mathbf{z}\|} = \mathbf{v}$$

exists and is nontrivial. Now let $s_{\nabla f}(\mathbf{y}, \mathbf{x})$ be a slope function for $\nabla f(\mathbf{x})$. Taking limits in

$$\begin{aligned} \mathbf{0} &= \frac{1}{\|\mathbf{z}_{m_k} - \mathbf{z}\|} [\nabla f(\mathbf{z}_{m_k}) - \nabla f(\mathbf{z})] \\ &= \frac{1}{\|\mathbf{z}_{m_k} - \mathbf{z}\|} s_{\nabla f}(\mathbf{z}_{m_k}, \mathbf{z})(\mathbf{z}_{m_k} - \mathbf{z}) \end{aligned}$$

then produces $\mathbf{0} = d^2 f(\mathbf{z})\mathbf{v}$. In other words, the second differential at \mathbf{z} is singular. If one can rule out such degeneracies, then all stationary points are isolated. Interested readers can consult the literature on Morse functions for further commentary on this subject [115].

12.5 Global Convergence of Block Relaxation

Verification of global convergence of block relaxation parallels the MM algorithm case. Careful scrutiny of the proof of Proposition 12.4.4 shows that it relies on five properties of the objective function $f(\mathbf{x})$ and the iteration map $M(\mathbf{x})$:

- (a) $f(\mathbf{x})$ is coercive on its convex open domain U ,
- (b) $f(\mathbf{x})$ has only isolated stationary points,
- (c) $M(\mathbf{x})$ is continuous,
- (d) \mathbf{y} is a fixed point of $M(\mathbf{x})$ if and only if it is a stationary point of $f(\mathbf{x})$,

(e) $f[M(\mathbf{y})] \leq f(\mathbf{y})$, with equality if and only if \mathbf{y} is a fixed point of $M(\mathbf{x})$.

Let us suppose for notational simplicity that the argument $\mathbf{x} = (\mathbf{v}, \mathbf{w})$ breaks into just two blocks. Criteria (a) and (b) can be demonstrated for many objective functions and are independent of the algorithm chosen to minimize $f(\mathbf{x})$. In block relaxation we ordinarily take U to be the Cartesian product $V \times W$ of two convex open sets. If we assume that $f(\mathbf{v}, \mathbf{w})$ is strictly convex in \mathbf{v} for fixed \mathbf{w} and vice versa, then the block relaxation updates are well defined. If $f(\mathbf{v}, \mathbf{w})$ is twice continuously differentiable, and $d_{\mathbf{v}}^2 f(\mathbf{v}, \mathbf{w})$ and $d_{\mathbf{w}}^2 f(\mathbf{v}, \mathbf{w})$ are invertible matrices, then application of the implicit function theorem demonstrates that the iteration map $M(\mathbf{x})$ is a composition of two differentiable maps. Criterion (c) is therefore valid. A fixed point $\mathbf{x} = (\mathbf{v}, \mathbf{w})$ satisfies the two equations $\nabla_{\mathbf{v}} f(\mathbf{v}, \mathbf{w}) = \mathbf{0}$ and $\nabla_{\mathbf{w}} f(\mathbf{v}, \mathbf{w}) = \mathbf{0}$, and criterion (d) follows. Finally, both block updates decrease $f(\mathbf{x})$. They give a strict decrease if and only if they actually change either argument \mathbf{v} or \mathbf{w} . Hence, criterion (e) is true. We emphasize that collectively these are sufficient but not necessary conditions. Observe that we have not assumed that $f(\mathbf{v}, \mathbf{w})$ is convex in both variables simultaneously.

12.6 Global Convergence of Gradient Algorithms

We now turn to the question of global convergence for gradient algorithms of the sort

$$\mathbf{x}_{m+1} = \mathbf{x}_m - A(\mathbf{x}_m)^{-1} \nabla f(\mathbf{x}_m).$$

The assumptions concerning $f(\mathbf{x})$ made in the previous section remain in force. A major impediment to establishing the global convergence of any minimization algorithm is the possible failure of the descent property

$$f(\mathbf{x}_{m+1}) \leq f(\mathbf{x}_m)$$

enjoyed by the MM algorithm. Provided the matrix $A(\mathbf{x}_m)$ is positive definite, the direction $\mathbf{v}_m = -A(\mathbf{x}_m)^{-1} \nabla f(\mathbf{x}_m)$ is guaranteed to point locally downhill. Hence, if we elect the natural strategy of instituting a limited line search along the direction \mathbf{v}_m emanating from \mathbf{x}_m , then we can certainly find an \mathbf{x}_{m+1} that decreases $f(\mathbf{x})$.

Although an exact line search is tempting, we may pay too great a price for precision when we merely need progress. The step-halving tactic mentioned in Chap. 10 is better than a full line search but not quite adequate for theoretical purposes. Instead, we require a sufficient decrease along a descent direction \mathbf{v} . This is summarized by the Armijo rule of considering only steps $t\mathbf{v}$ satisfying the inequality

$$f(\mathbf{x} + t\mathbf{v}) \leq f(\mathbf{x}) + \alpha t df(\mathbf{x})\mathbf{v} \quad (12.10)$$

for $t > 0$ and some fixed α in $(0, 1)$. To avoid too stringent a test, we take a low value of α such as 0.01. In combining Armijo's rule with regular step decrementing, we first test the step \mathbf{v} . If it satisfies Armijo's rule we are done. If it fails, we choose $\sigma \in (0, 1)$ and test $\sigma\mathbf{v}$. If this fails, we test $\sigma^2\mathbf{v}$, and so forth until we encounter and take the first partial step $\sigma^k\mathbf{v}$ that works. In step halving, obviously $\sigma = 1/2$.

Step halving can be combined with a partial line search. For instance, suppose the line search has been confined to the interval $t \in [0, s]$. If the point $\mathbf{x} + s\mathbf{v}$ passes Armijo's test, then we accept it. Otherwise, we fit a cubic to the function $t \mapsto f(\mathbf{x} + t\mathbf{v})$ on the interval $[0, s]$ as described in Sect. 11.4. If the minimum point t of the cubic approximation satisfies $t \geq \sigma s$ and passes Armijo's test, then we accept $\mathbf{x} + t\mathbf{v}$. Otherwise, we replace the interval $[0, s]$ by the interval $[0, \sigma s]$ and proceed inductively. For the sake of simplicity in the sequel, we will ignore this elaboration of step halving and concentrate on the unadorned version.

We would like some guarantee that the exponent k of the step decrementing power σ^k does not grow too large. Mindful of this criterion, we suppose that the positive definite matrix $A(\mathbf{x})$ depends continuously on \mathbf{x} . This is not much of a restriction for Newton's method, the Gauss-Newton algorithm, the MM gradient algorithm, or scoring. If we combine continuity with coerciveness, then we can conclude that there exist positive constants β , γ , δ , and ϵ with

$$\begin{aligned} \|A(\mathbf{x})\| &\leq \beta, & \|A(\mathbf{x})^{-1}\| &\leq \gamma \\ \|\nabla f(\mathbf{x})\| &\leq \epsilon, & \|s_f^2(\mathbf{y}, \mathbf{x})\| &\leq \delta \end{aligned}$$

for all \mathbf{x} and \mathbf{y} in the compact set $D = \{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$ where any descent algorithm acts. Here $s_f^2(\mathbf{y}, \mathbf{x})$ is the second slope of $f(\mathbf{x})$.

Before we tackle Armijo's rule, let us consider the more pressing question of whether the proposed points $\mathbf{x} + \mathbf{v}$ lie in the domain U of $f(\mathbf{x})$. This is too much to hope for, but it is worth considering whether $\mathbf{x} + \sigma^d\mathbf{v}$ always lies in U for some fixed power σ^d . Fortunately, $\mathbf{v}(\mathbf{x}) = -A(\mathbf{x})^{-1}\nabla f(\mathbf{x})$ satisfies the bound

$$\|\mathbf{v}(\mathbf{x})\| \leq \gamma\epsilon$$

on D . Now suppose no single power σ^k is adequate for all $\mathbf{x} \in D$. Then there exists a sequence of points $\mathbf{x}_k \in D$ with $\mathbf{y}_k = \mathbf{x}_k + \sigma^k\mathbf{v}(\mathbf{x}_k) \notin U$. Passing to a subsequence if necessary, we can assume that \mathbf{x}_k converges to $\mathbf{x} \in D$. Because σ^k is tending to 0, and $\mathbf{v}(\mathbf{x})$ is bounded on D , the sequence \mathbf{y}_k likewise converges to \mathbf{x} . Since the complement of U is closed, \mathbf{x} must lie in the complement of U as well as in D . This contradiction proves our contention.

To use these bounds, let $\mathbf{v} = -A(\mathbf{x})^{-1}\nabla f(\mathbf{x})$ and consider the inequality

$$f(\mathbf{x} + t\mathbf{v}) = f(\mathbf{x}) + tdf(\mathbf{x})\mathbf{v} + \frac{1}{2}t^2\mathbf{v}^*s_f^2(\mathbf{x} + t\mathbf{v}, \mathbf{x})\mathbf{v}$$

$$\leq f(\mathbf{x}) + tdf(\mathbf{x})\mathbf{v} + \frac{1}{2}t^2\delta\|\mathbf{v}\|^2 \tag{12.11}$$

for \mathbf{x} and $\mathbf{x} + t\mathbf{v}$ in D . Taking into account the bound on $\|A(\mathbf{x})\|$ and the identity

$$\|A(\mathbf{x})^{1/2}\| = \|A(\mathbf{x})\|^{1/2}$$

entailed by Proposition 2.2.1, we also have

$$\begin{aligned} \|\nabla f(\mathbf{x})\|^2 &= \|A(\mathbf{x})^{1/2}A(\mathbf{x})^{-1/2}\nabla f(\mathbf{x})\|^2 \\ &\leq \|A(\mathbf{x})^{1/2}\|^2\|A(\mathbf{x})^{-1/2}\nabla f(\mathbf{x})\|^2 \\ &\leq \beta df(\mathbf{x})A(\mathbf{x})^{-1}\nabla f(\mathbf{x}). \end{aligned} \tag{12.12}$$

It follows that

$$\begin{aligned} \|\mathbf{v}\|^2 &= \|A(\mathbf{x})^{-1}\nabla f(\mathbf{x})\|^2 \\ &\leq \gamma^2\|\nabla f(\mathbf{x})\|^2 \\ &\leq -\beta\gamma^2df(\mathbf{x})\mathbf{v}. \end{aligned}$$

Combining this last inequality with inequality (12.11) yields

$$f(\mathbf{x} + t\mathbf{v}) \leq f(\mathbf{x}) + t\left(1 - \frac{\beta\gamma^2\delta}{2}t\right)df(\mathbf{x})\mathbf{v}.$$

Hence, as soon as σ^k satisfies

$$1 - \frac{\beta\gamma^2\delta}{2}\sigma^k \geq \alpha,$$

Armijo’s rule (12.10) holds. In terms of k , backtracking is guaranteed to succeed in at most

$$k_{\max} = \max\left\{\left\lceil \frac{1}{\ln \sigma} \ln \frac{2(1-\alpha)}{\beta\gamma^2\delta} \right\rceil, d\right\}$$

decrements. Of course, a lower value of k may suffice.

Proposition 12.6.1 *Suppose that all stationary points of $f(\mathbf{x})$ are isolated and that the stated continuity, differentiability, positive definiteness, and coerciveness assumptions are true. Then any sequence of iterates \mathbf{x}_m generated by the iteration map $M(\mathbf{x}) = \mathbf{x} - tA(\mathbf{x})^{-1}\nabla f(\mathbf{x})$ with t chosen by step decrementing possesses a limit, and that limit is a stationary point of $f(\mathbf{x})$. If $f(\mathbf{x})$ is strictly convex, then $\lim_{m \rightarrow \infty} \mathbf{x}_m$ is the minimum point.*

Proof: Let $\mathbf{v}_m = -A(\mathbf{x}_m)^{-1}\nabla f(\mathbf{x}_m)$ and $\mathbf{x}_{m+1} = \mathbf{x}_m + \sigma^{k_m}\mathbf{v}_m$. The sequence $f(\mathbf{x}_m)$ is decreasing by construction. Because the function $f(\mathbf{x})$ is bounded below on the compact set $D = \{\mathbf{x} \in U : f(\mathbf{x}) \leq f(\mathbf{x}_0)\}$,

$f(\mathbf{x}_m)$ is bounded below as well and possesses a limit. Based on Armijo's rule (12.10) and inequality (12.12), we calculate

$$\begin{aligned} f(\mathbf{x}_m) - f(\mathbf{x}_{m+1}) &\geq -\alpha\sigma^{k_m} df(\mathbf{x}_m)\mathbf{v}_m \\ &= \alpha\sigma^{k_m} df(\mathbf{x}_m)A(\mathbf{x}_m)^{-1}\nabla f(\mathbf{x}_m) \\ &\geq \frac{\alpha\sigma^{k_m}}{\beta} \|\nabla f(\mathbf{x}_m)\|^2. \end{aligned}$$

Since $\sigma^{k_m} \geq \sigma^{k_{\max}}$, and the difference $f(\mathbf{x}_m) - f(\mathbf{x}_{m+1})$ tends to 0, we deduce that $\|\nabla f(\mathbf{x}_m)\|$ tends to 0. This conclusion and the inequality

$$\begin{aligned} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| &= \sigma^{k_m} \|A(\mathbf{x}_m)^{-1}\nabla f(\mathbf{x}_m)\| \\ &\leq \sigma^{k_m} \gamma \|\nabla f(\mathbf{x}_m)\|, \end{aligned}$$

demonstrate that $\|\mathbf{x}_{m+1} - \mathbf{x}_m\|$ tends to 0 as well. Given these results, Propositions 12.4.2 and 12.4.3 are true. All claims of the current proposition now follow as in the proof of Proposition 12.4.4. ■

12.7 Problems

1. Consider the functions $f(x) = x - x^3$ and $g(x) = x + x^3$ on \mathbb{R} . Show that the iterates $x_{m+1} = f(x_m)$ are locally attracted to 0 and that the iterates $x_{m+1} = g(x_m)$ are locally repelled by 0. In both cases $f'(0) = g'(0) = 1$.
2. Consider the iteration map $h(x) = \sqrt{a+x}$ on $(0, \infty)$ for $a > 0$. Find the fixed point of $h(x)$ and show that it is locally attractive. Is it also globally attractive?
3. In Example 10.2.1 suppose $x_0 = 1$ and $a \in (0, 2)$. Demonstrate that

$$\begin{aligned} x_m &= \frac{1 - (1-a)^{2^m}}{a} \\ \left|x_{m+1} - \frac{1}{a}\right| &= a \left|x_m - \frac{1}{a}\right|^2. \end{aligned}$$

This shows very explicitly that x_m converges to $1/a$ at a quadratic rate.

4. In Example 10.2.2 prove that

$$\begin{aligned} x_m &= \sqrt{a} + \frac{2\sqrt{a}}{\left[\left(1 + \frac{2\sqrt{a}}{x_0 - \sqrt{a}}\right)^{2^m} - 1\right]} \\ \left|x_{m+1} - \sqrt{a}\right| &\leq \frac{1}{2\sqrt{a}} \left|x_m - \sqrt{a}\right|^2 \end{aligned}$$

when $n = 2$ and $x_0 > 0$. Thus, Newton's method converges at a quadratic rate. Use the first of these formulas or the iteration equation directly to show that $\lim_{m \rightarrow \infty} x_m = -\sqrt{a}$ for $x_0 < 0$.

5. Suppose the real-valued function $f(x)$ is twice continuously differentiable on the interval (a, b) with a root y where $f(y) = 0$ and $f'(y) \neq 0$. Show that the iteration scheme

$$x_{n+1} = x_n - \frac{f(x_n)^2}{f[x_n + f(x_n)] - f(x_n)}$$

converges at a quadratic rate to y if x_0 is sufficiently close to y . In particular, demonstrate that

$$\lim_{n \rightarrow \infty} \frac{x_{n+1} - y}{(x_n - y)^2} = \frac{f''(y)[f'(y) - 1]}{2f'(y)}.$$

6. Let \mathbf{A} and \mathbf{B} be $n \times n$ matrices. If \mathbf{A} and $\mathbf{A} - \mathbf{B}^* \mathbf{A} \mathbf{B}$ are both positive definite, then show that \mathbf{B} has spectral radius $\rho(\mathbf{B}) < 1$. Note that \mathbf{A} is symmetric, but \mathbf{B} need not be symmetric. (Hint: Consider the quadratic form $\mathbf{v}^*(\mathbf{A} - \mathbf{B}^* \mathbf{A} \mathbf{B})\mathbf{v}$ for an eigenvector \mathbf{v} of \mathbf{B} .)
7. In block relaxation with b blocks, let $B_i(\mathbf{x})$ be the map that updates block i and leaves the other blocks fixed. Show that the overall iteration map $M(\mathbf{x}) = B_b \circ \cdots \circ B_1(\mathbf{x})$ has differential $dB_b(\mathbf{y}) \cdots dB_1(\mathbf{y})$ at a fixed point \mathbf{y} . Write $dB_i(\mathbf{y})$ as a block matrix and identify the blocks by applying the implicit function theorem as needed. Do not confuse $B_i(\mathbf{x})$ with the update $M_i(\mathbf{x})$ of the text. In fact, $M_i(\mathbf{x})$ only summarizes the update of block i , and its argument is the value of \mathbf{x} at the start of the current round of updates.
8. Consider a Poisson-distributed random variable Y with mean $a\theta + b$, where a and b are known positive constants and $\theta \geq 0$ is a parameter to be estimated. An EM algorithm for estimating θ can be concocted that takes as complete data independent Poisson random variables U and V with means $a\theta$ and b and sum $U + V = Y$. If $Y = y$ is observed, then show that the EM iterates are defined by

$$\theta_{m+1} = \frac{y\theta_m}{a\theta_m + b}.$$

Show that these iterates converge monotonically to the maximum likelihood estimate $\max\{0, (y - b)/a\}$. When $y = b$, verify that convergence to the boundary value 0 occurs at a rate slower than linear [90]. (Hint: When $y = b$, check that $\theta_{m+1} = b\theta_0/(ma\theta_0 + b)$.)

9. The sublinear convergence of the EM algorithm exhibited in the previous problem occurs in other problems. Here is a conceptually harder

example by Robert Jennrich. Suppose that W_1, \dots, W_n and B are independent normally distributed random variables with 0 means. Let σ_w^2 be the common variance of the W_i and σ_b^2 be the variance of B . If the values y_i of the linear combinations $Y_i = B + W_i$ are observed, then show that the EM algorithm amounts to

$$\begin{aligned} \sigma_{m+1,b}^2 &= \left(\frac{n\sigma_{mb}^2 \bar{y}}{n\sigma_{mb}^2 + \sigma_{mw}^2} \right)^2 + \frac{\sigma_{mb}^2 \sigma_{nw}^2}{n\sigma_{mb}^2 + \sigma_{mw}^2} \\ \sigma_{m+1,w}^2 &= \frac{n-1}{n} s_y^2 + \left(\frac{\sigma_{mw}^2 \bar{y}}{n\sigma_{mb}^2 + \sigma_{mw}^2} \right)^2 + \frac{\sigma_{mb}^2 \sigma_{mw}^2}{n\sigma_{mb}^2 + \sigma_{mw}^2}, \end{aligned}$$

where $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ and $s_y^2 = \frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2$ are the sample mean and variance. Although one can formally calculate the maximum likelihood estimates $\hat{\sigma}_w^2 = s_y^2$ and $\hat{\sigma}_b^2 = \bar{y}^2 - s_y^2/n$, these are only valid provided $\hat{\sigma}_b^2 \geq 0$. If for instance $\bar{y} = 0$, then the EM iterates will converge to $\sigma_w^2 = (n-1)s_y^2/n$ and $\sigma_b^2 = 0$. Show that convergence is sublinear when $\bar{y} = 0$.

10. Consider independent observations y_1, \dots, y_n from the univariate t -distribution. These data have loglikelihood

$$\begin{aligned} L &= -\frac{n}{2} \ln \sigma^2 - \frac{\nu+1}{2} \sum_{i=1}^n \ln(\nu + \delta_i^2) \\ \delta_i^2 &= \frac{(y_i - \mu)^2}{\sigma^2}. \end{aligned}$$

To illustrate the occasionally bizarre behavior of the MM algorithm, we take $\nu = 0.05$ and the data vector $y = (-20, 1, 2, 3)^*$ with $n = 4$ observations. Devise an MM maximum likelihood algorithm for estimating μ with σ^2 fixed at 1. Show that the iteration map is

$$\begin{aligned} \mu_{m+1} &= \frac{\sum_{i=1}^n w_{mi} y_i}{\sum_{i=1}^n w_{mi}} \\ w_{mi} &= \frac{\nu + 1}{\nu + (y_i - \mu_m)^2}. \end{aligned}$$

Plot the likelihood curve and show that it has the four local maxima $-19.993, 1.086, 1.997,$ and 2.906 and the three local minima $-14.516, 1.373,$ and 2.647 . Demonstrate numerically convergence to a local maximum that is not the global maximum. Show that the algorithm converges to a local minimum in one step starting from -1.874 or -0.330 [191].

11. Suppose the data displayed in Table 12.1 constitute a random sample from a bivariate normal distribution with both means 0, variances σ_1^2

TABLE 12.1. Bivariate normal data for the EM algorithm

Obs	Obs	Obs	Obs	Obs	Obs
(1,1)	(1, -1)	(-1, 1)	(-1, -1)	(2,*)	(2,*)
(-2,*)	(-2,*)	(*,2)	(*,2)	(*, -2)	(*, -2)

and σ_2^2 , and correlation coefficient ρ . The asterisks indicate missing values. Specify the EM algorithm for estimating σ_1^2 , σ_2^2 , and ρ . Show that the observed loglikelihood has a saddle point at the point where $\rho = 0$ and $\sigma_1^2 = \sigma_2^2 = \frac{5}{2}$. If the EM algorithm starts with $\rho = 0$, prove that convergence to the saddle point occurs [275]. (Hints: At the given point, show that the observed loglikelihood achieves a global maximum in σ_1^2 and σ_2^2 with ρ fixed at 0 and a local minimum in ρ with σ_1^2 and σ_2^2 fixed at $\frac{5}{2}$. Also show that the EM algorithm iterates satisfy

$$\begin{aligned} \sigma_{m+1,1}^2 - \frac{5}{2} &= \frac{1}{3} \left(\sigma_{m1}^2 - \frac{5}{2} \right) \\ \sigma_{m+1,2}^2 - \frac{5}{2} &= \frac{1}{3} \left(\sigma_{m2}^2 - \frac{5}{2} \right) \\ \rho_{m+1} &= 0 \end{aligned}$$

along the slice $\rho = 0$.)

12. Under the hypotheses of Proposition 12.4.4, if the MM gradient algorithm is started close enough to a local minimum \mathbf{y} of $f(\mathbf{x})$, then the iterates \mathbf{x}_m converge to \mathbf{y} without step decrementing. Prove that for all sufficiently large m , either $\mathbf{x}_m = \mathbf{y}$ or $f(\mathbf{x}_{m+1}) < f(\mathbf{x}_m)$ [163]. (Hints: Let $\mathbf{v}_m = \mathbf{x}_{m+1} - \mathbf{x}_m$, $\mathbf{C}_m = s_f^2(\mathbf{x}_{m+1}, \mathbf{x}_m)$, and $\mathbf{D}_m = d^2g(\mathbf{x}_m | \mathbf{x}_m)$. Show that

$$f(\mathbf{x}_{m+1}) = f(\mathbf{x}_m) + \frac{1}{2} \mathbf{v}_m^* (\mathbf{C}_m - 2\mathbf{D}_m) \mathbf{v}_m.$$

Then use a continuity argument, noting that $d^2g(\mathbf{y} | \mathbf{y}) - d^2f(\mathbf{y})$ is positive semidefinite and $d^2g(\mathbf{y} | \mathbf{y})$ is positive definite.)

13. Let $M(\mathbf{x})$ be the MM algorithm or MM gradient algorithm map. Consider the modified algorithm $M_t(\mathbf{x}) = \mathbf{x} + t[M(\mathbf{x}) - \mathbf{x}]$ for $t > 0$. At a local optimum \mathbf{y} , show that the spectral radius ρ_t of the differential $dM_t(\mathbf{y}) = (1 - t)\mathbf{I} + tdM(\mathbf{y})$ satisfies $\rho_t < 1$ when $0 < t < 2$. Hence, Ostrowski's theorem implies local attraction of $M_t(\mathbf{x})$ to \mathbf{y} . If the largest and smallest eigenvalues of $dM(\mathbf{y})$ are ω_{\max} and ω_{\min} , then prove that ρ_t is minimized by taking $t = [1 - (\omega_{\min} + \omega_{\max})/2]^{-1}$. In practice, the eigenvalues of $dM(\mathbf{y})$ are impossible to predict without advance knowledge of \mathbf{y} , but for many problems the value $t = 2$

works well [163]. (Hints: All eigenvalues of $dM(\mathbf{y})$ occur on $[0, 1)$. To every eigenvalue ω of $dM(\mathbf{y})$, there corresponds an eigenvalue $\omega_t = 1 - t + t\omega$ of $dM_t(\mathbf{y})$ and vice versa.)

14. In the notation of Chap. 9, prove the EM algorithm formula

$$d^2L(\boldsymbol{\theta}) = d^{20}Q(\boldsymbol{\theta} \mid \boldsymbol{\theta}) + \text{Var}[\nabla \ln f(X \mid \boldsymbol{\theta}) \mid Y, \boldsymbol{\theta}]$$

of Louis [180].

15. Which of the following functions is coercive on its domain?

- (a) $f(x) = x + 1/x$ on $(0, \infty)$,
- (b) $f(x) = x - \ln x$ on $(0, \infty)$,
- (c) $f(\mathbf{x}) = x_1^2 + x_2^2 - 2x_1x_2$ on \mathbb{R}^2 ,
- (d) $f(\mathbf{x}) = x_1^4 + x_2^4 - 3x_1x_2$ on \mathbb{R}^2 ,
- (e) $f(\mathbf{x}) = x_1^2 + x_2^2 + x_3^2 - \sin(x_1x_2x_3)$ on \mathbb{R}^3 .

Give convincing reasons in each case.

16. Consider a polynomial $p(\mathbf{x})$ in n variables x_1, \dots, x_n . Suppose that $p(\mathbf{x}) = \sum_{i=1}^n c_i x_i^{2m} + \text{lower-order terms}$, where all $c_i > 0$ and where a lower-order term is a product $b x_1^{m_1} \cdots x_n^{m_n}$ with $\sum_{i=1}^n m_i < 2m$. Prove rigorously that $p(\mathbf{x})$ is coercive on \mathbb{R}^n .
17. Demonstrate that $h(\mathbf{x}) + k(\mathbf{x})$ is coercive on \mathbb{R}^n if $k(\mathbf{x})$ is convex and $h(\mathbf{x})$ satisfies $\lim_{\|\mathbf{x}\| \rightarrow \infty} \|\mathbf{x}\|^{-1} h(\mathbf{x}) = \infty$. Problem 31 of Chap. 6 is a special case. (Hint: Apply definition (6.4) of Chap. 6.)
18. In some problems it is helpful to broaden the notion of coerciveness. Consider a continuous function $f : \mathbb{R}^n \mapsto \mathbb{R}$ such that the limit

$$c = \lim_{r \rightarrow \infty} \inf_{\|\mathbf{x}\| \geq r} f(\mathbf{x})$$

exists. The value of c can be finite or ∞ but not $-\infty$. Now let \mathbf{y} be any point with $f(\mathbf{y}) < c$. Show that the set $S_{\mathbf{y}} = \{\mathbf{x} : f(\mathbf{x}) \leq f(\mathbf{y})\}$ is compact and that $f(\mathbf{x})$ attains its global minimum on $S_{\mathbf{y}}$. The particular function

$$g(\mathbf{x}) = \frac{x_1 + 2x_2}{1 + x_1^2 + x_2^2}$$

furnishes an example when $n = 2$. Demonstrate that the limit c equals 0. What is the minimum value and minimum point of $g(\mathbf{x})$? (Hint: What is the minimum value of $g(\mathbf{x})$ on the circle $\{\mathbf{x} : \|\mathbf{x}\| = r\}$?)

19. Let $f(\mathbf{x})$ be a convex function on \mathbb{R}^n . Prove that all sublevel sets $\{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) \leq b\}$ are bounded if and only if

$$\lim_{r \rightarrow \infty} \inf_{\|\mathbf{x}\| \geq r} f(\mathbf{x}) > 0.$$

20. Assume that the function $f(\mathbf{x})$ is (a) continuously differentiable, (b) maps \mathbb{R}^n into itself, (c) has Jacobian $df(\mathbf{x})$ of full rank at each \mathbf{x} , and (d) has coercive norm $\|f(\mathbf{x})\|$. Show that the image of $f(\mathbf{x})$ is all of \mathbb{R}^n . (Hints: The image is connected. Show that it is open via the inverse function theorem and closed because of coerciveness.)
21. Consider a sequence \mathbf{x}_m in \mathbb{R}^n . Verify that the set of cluster points of \mathbf{x}_m is closed. If \mathbf{x}_m is bounded, then show that it has a limit if and only if it has at most one cluster point.
22. In our exposition of least absolute deviation regression, we considered in Problem 14 of Chap. 8 a modified iteration scheme that minimizes the criterion

$$h_\epsilon(\boldsymbol{\theta}) = \sum_{i=1}^p \{[y_i - \mu_i(\boldsymbol{\theta})]^2 + \epsilon\}^{1/2}. \quad (12.13)$$

For a sequence of constants ϵ_m tending to 0, let $\boldsymbol{\theta}_m$ be a corresponding sequence minimizing (12.13). If $\boldsymbol{\phi}$ is a cluster point of this sequence and the regression functions $\mu_i(\boldsymbol{\theta})$ are continuous, then show that $\boldsymbol{\phi}$ minimizes $h_0(\boldsymbol{\theta}) = \sum_{i=1}^p |y_i - \mu_i(\boldsymbol{\theta})|$. If, in addition, the minimum point $\boldsymbol{\phi}$ of $h_0(\boldsymbol{\theta})$ is unique and $\lim_{\|\boldsymbol{\theta}\| \rightarrow \infty} \sum_{i=1}^p |\mu_i(\boldsymbol{\theta})| = \infty$, then prove that $\lim_{m \rightarrow \infty} \boldsymbol{\theta}_m = \boldsymbol{\phi}$. (Hints: For the first assertion, take limits in

$$h_\epsilon(\boldsymbol{\theta}_m) \leq h_\epsilon(\boldsymbol{\theta}).$$

For the second assertion, it suffices to prove that the sequence $\boldsymbol{\theta}_m$ is confined to a bounded set. To prove this fact, demonstrate the inequalities $|\mu| + |y| + \sqrt{\epsilon} \geq \sqrt{r^2 + \epsilon} \geq |\mu| - |y|$ for $r = y - \mu$.)

23. Example 10.2.5 and Problem 9 of Chap. 10 suggest a method of accelerating the MM gradient algorithm. Suppose we are maximizing the loglikelihood $L(\boldsymbol{\theta})$ using the surrogate function is $g(\boldsymbol{\theta} | \boldsymbol{\theta}_n)$. To accelerate the MM gradient algorithm, we can replace the positive definite matrix $B(\boldsymbol{\theta})^{-1} = -d^{20}g(\boldsymbol{\theta} | \boldsymbol{\theta})$ by a matrix that better approximates the observed information $A(\boldsymbol{\theta}) = -d^2L(\boldsymbol{\theta})$. Note that often $d^{20}g(\boldsymbol{\theta} | \boldsymbol{\theta})$ is diagonal and therefore trivial to invert. Now consider the formal expansion

$$\mathbf{A}^{-1} = (\mathbf{B}^{-1} + \mathbf{A} - \mathbf{B}^{-1})^{-1}$$

$$\begin{aligned}
&= \{ \mathbf{B}^{-\frac{1}{2}} [\mathbf{I} - \mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{-1} - \mathbf{A}) \mathbf{B}^{\frac{1}{2}}] \mathbf{B}^{-\frac{1}{2}} \}^{-1} \\
&= \mathbf{B}^{\frac{1}{2}} \sum_{i=0}^{\infty} [\mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{-1} - \mathbf{A}) \mathbf{B}^{\frac{1}{2}}]^i \mathbf{B}^{\frac{1}{2}}.
\end{aligned}$$

If we truncate this series after a finite number of terms, then we recover the first iterate of equation (10.18) in the disguised form

$$\mathbf{S}_j = \mathbf{B}^{\frac{1}{2}} \sum_{i=0}^j [\mathbf{B}^{\frac{1}{2}} (\mathbf{B}^{-1} - \mathbf{A}) \mathbf{B}^{\frac{1}{2}}]^i \mathbf{B}^{\frac{1}{2}}.$$

The accelerated algorithm

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + \mathbf{S}_j(\boldsymbol{\theta}_m) \nabla L(\boldsymbol{\theta}_m) \quad (12.14)$$

has several desirable properties.

- Show that \mathbf{S}_j is positive definite and hence that the update (12.14) is an ascent algorithm. (Hint: Use the fact that $\mathbf{B}^{-1} - \mathbf{A}$ is positive semidefinite.)
- Algorithm (12.14) has differential

$$\mathbf{I} + \mathbf{S}_j(\boldsymbol{\theta}_\infty) d^2 L(\boldsymbol{\theta}_\infty) = \mathbf{I} - \mathbf{S}_j(\boldsymbol{\theta}_\infty) \mathbf{A}(\boldsymbol{\theta}_\infty)$$

at a local maximum $\boldsymbol{\theta}_\infty$. If $d^2 L(\boldsymbol{\theta}_\infty)$ is negative definite, then prove that all eigenvalues of this differential lie on $[0, 1)$. (Hint: The eigenvalues are determined by the stationary points of the Rayleigh quotient $\mathbf{v}^* [\mathbf{A}^{-1}(\boldsymbol{\theta}_\infty) - \mathbf{S}_j(\boldsymbol{\theta}_\infty)] \mathbf{v} / \mathbf{v}^* \mathbf{A}^{-1}(\boldsymbol{\theta}_\infty) \mathbf{v}$.)

- If ρ_j is the spectral radius of the differential, then demonstrate that $\rho_j \leq \rho_{j-1}$, with strict inequality when $\mathbf{B}^{-1}(\boldsymbol{\theta}_\infty) - \mathbf{A}(\boldsymbol{\theta}_\infty)$ is positive definite.

In other words, the accelerated algorithm (12.14) is guaranteed to converge faster than the MM gradient algorithm. It will be particularly useful for maximum likelihood problems with many parameters because it entails no matrix inversion or multiplication, just matrix times vector multiplication. When $j = 1$, it takes the simple form

$$\boldsymbol{\theta}_{m+1} = \boldsymbol{\theta}_m + [2\mathbf{B}(\boldsymbol{\theta}_m) - \mathbf{B}(\boldsymbol{\theta}_m) \mathbf{A}(\boldsymbol{\theta}_m) \mathbf{B}(\boldsymbol{\theta}_m)] \nabla L(\boldsymbol{\theta}_m).$$

13

Penalty and Barrier Methods

13.1 Introduction

Penalties and barriers feature prominently in two areas of modern optimization theory. First, both devices are employed to solve constrained optimization problems [96, 183, 226]. The general idea is to replace hard constraints by penalties or barriers and then exploit the well-oiled machinery for solving unconstrained problems. Penalty methods operate on the exterior of the feasible region and barrier methods on the interior. The strength of a penalty or barrier is determined by a tuning constant. In classical penalty methods, a single global tuning constant is gradually sent to ∞ ; in barrier methods, it is gradually sent to 0. Nothing prevents one from assigning different tuning constants to different penalties or barriers in the same problem. Either strategy generates a sequence of solutions that converges in practice to the solution of the original constrained optimization problem.

One of the lessons of the current chapter is that it is profitable to view penalties and barriers from the perspective of the MM algorithm. For example, this mental exercise suggests a way of engineering barrier tuning constants in a constrained minimization problem so that the objective function is forced steadily downhill [41, 162, 253]. Over time the tuning constant for each inequality constraint adapts to the need to avoid the constraint or converge to it.

A detailed study of penalty methods in constrained optimization requires considerable knowledge of the convex calculus. For this reason we defer

to Chap. 16 our presentation of exact penalty methods. These methods substitute absolute value and hinge penalties for square penalties. Numerical analysts have shied away from exact penalty methods because of the difficulties in working with non-differentiable functions. Such pessimism is unwarranted, however. In convex programs one can easily follow the solution path as the penalty constant increases. Certainly, path following has been highly successful in interior point programming. Readers who want to pursue this traditional application of path following should consult one or more of the superb references [19, 30, 96, 203, 205, 274].

Imposition of penalties has other beneficial effects. The addition of convex penalties can regularize a problem by reducing the effect of noise and eliminating spurious local minima. Penalties also steer solutions to unconstrained optimization problems in productive directions. Recall, for instance, the beneficial effects of smoothing penalties in transmission tomography. Bayesian statistics is predicated on the philosophy that prior knowledge should never be ignored. Beyond smoothing and exploitation of prior knowledge, penalties play a role in model selection. In lasso penalized estimation, statisticians impose an ℓ_1 penalty that shrinks parameter estimates toward zero and performs a kind of continuous model selection [75, 256]. The predictors whose estimated regression coefficients are exactly zero are candidates for elimination from the model. With the enormous data sets now confronting statisticians, considerations of model parsimony have taken on greater urgency. In addition to this philosophical justification, imposition of lasso penalties also has an huge impact on computational speed. Standard methods of regression require matrix diagonalization, matrix inversion, or, at the very least, the solution of large systems of linear equations. Because the number of arithmetic operations for these processes scales as the cube of the number of predictors, problems with tens of thousands of predictors appear intractable. Recent research has shown this assessment to be too pessimistic [36, 83, 119, 143, 210, 267]. Coordinate descent methods mesh well with the lasso and are simple, fast, and stable. We will see how their potential to transform data mining plays out in both ℓ_1 and ℓ_2 regression.

13.2 Rudiments of Barrier and Penalty Methods

In general, unconstrained optimization problems are easier to solve than constrained optimization problems, and equality constrained problems are easier to solve than inequality constrained problems. To simplify analysis, mathematical scientists rely on several devices. For instance, one can replace the inequality constraint $g(\mathbf{x}) \leq 0$ by the equality constraint $g_+(\mathbf{x}) = 0$, where $g_+(\mathbf{x}) = \max\{g(\mathbf{x}), 0\}$. This tactic is not entirely satisfactory because $g_+(\mathbf{x})$ has kinks along the boundary $g(\mathbf{x}) = 0$. The smoother substitute $g_+(\mathbf{x})^2$ avoids the kinks in first derivatives. Alternatively, one

can introduce an extra parameter y and require $g(\mathbf{x}) + y = 0$ and $y \geq 0$. This tactic substitutes a simple inequality constraint for a complex inequality constraint.

The addition of barrier and penalty terms to the objective function $f(\mathbf{x})$ is a more systematic approach. Later in the chapter we will discuss the role of penalties in producing sparse solutions. In the current section, penalties are introduced to steer the optimization process toward the feasible region. In the penalty method we construct a continuous nonnegative penalty $p(\mathbf{x})$ that is 0 on the feasible region and positive outside it. We then optimize the functions $f(\mathbf{x}) + \lambda_m p(\mathbf{x})$ for an increasing sequence of tuning constants λ_m that tend to ∞ . The penalty method works from the outside of the feasible region inward. Under the right hypotheses, the sequence of unconstrained solutions \mathbf{x}_m tends to a solution of the constrained optimization problem.

Example 13.2.1 *Linear Regression with Linear Constraints*

Consider the problem of minimizing $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ subject to the linear constraints $\mathbf{V}\boldsymbol{\beta} = \mathbf{d}$. If we take the penalty function $p(\boldsymbol{\beta}) = \|\mathbf{V}\boldsymbol{\beta} - \mathbf{d}\|^2$, then we must minimize at each stage the function

$$h_m(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda_m \|\mathbf{V}\boldsymbol{\beta} - \mathbf{d}\|^2.$$

Setting the gradient

$$\nabla h_m(\boldsymbol{\beta}) = -2\mathbf{X}^*(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) + 2\lambda_m \mathbf{V}^*(\mathbf{V}\boldsymbol{\beta} - \mathbf{d})$$

equal to $\mathbf{0}$ yields the sequence of solutions

$$\boldsymbol{\beta}_m = (\mathbf{X}^*\mathbf{X} + \lambda_m \mathbf{V}^*\mathbf{V})^{-1}(\mathbf{X}^*\mathbf{y} + \lambda_m \mathbf{V}^*\mathbf{d}). \quad (13.1)$$

In a moment we will demonstrate that the $\boldsymbol{\beta}_m$ tend to the constrained solution as λ_m tends to ∞ . ■

In contrast, the barrier method works from the inside of the feasible region outward by introducing a continuous barrier function $b(\mathbf{x})$ that is finite on the interior of the feasible region and infinite on its boundary. We then optimize the sequence of functions $f(\mathbf{x}) + \mu_m b(\mathbf{x})$ as the decreasing sequence of tuning constants μ_m tends to 0. Again under the right hypotheses, the sequence of unconstrained solutions \mathbf{x}_m tends to the solution of the constrained optimization problem.

Example 13.2.2 *Estimation of Multinomial Proportions*

In estimating multinomial proportions, we minimize the negative loglikelihood $-\sum_{i=1}^d n_i \ln p_i$ subject to the constraints $\sum_{i=1}^d p_i = 1$ and $p_i \geq 0$ for d categories. An appropriate barrier function is $-\sum_{i=1}^d \ln p_i$. The minimum of the function

$$h_m(\mathbf{p}) = -\sum_{i=1}^d n_i \ln p_i - \mu_m \sum_{i=1}^d \ln p_i$$

subject to the constraint $\sum_{i=1}^d p_i = 1$ occurs at the point with coordinates

$$p_{mi} = \frac{n_i + \mu_m}{n + d\mu_m},$$

where $n = \sum_{i=1}^d n_i$. In this example, it is clear that the solution vector \mathbf{p}_m occurs on the interior of the parameter space and tends to the maximum likelihood estimate as μ_m tends to 0. ■

The next proposition highlights the ascent and descent properties of the penalty and barrier methods.

Proposition 13.2.1 *Consider two real-valued functions $f(\mathbf{x})$ and $g(\mathbf{x})$ on a common domain and two positive constants $\alpha < \beta$. Suppose the linear combination $f(\mathbf{x}) + \alpha g(\mathbf{x})$ attains its minimum value at \mathbf{y} and the linear combination $f(\mathbf{x}) + \beta g(\mathbf{x})$ attains its minimum value at \mathbf{z} . Then we have $f(\mathbf{y}) \leq f(\mathbf{z})$ and $g(\mathbf{y}) \geq g(\mathbf{z})$.*

Proof: Adding the two inequalities

$$\begin{aligned} f(\mathbf{z}) + \beta g(\mathbf{z}) &\leq f(\mathbf{y}) + \beta g(\mathbf{y}) \\ -f(\mathbf{z}) - \alpha g(\mathbf{z}) &\leq -f(\mathbf{y}) - \alpha g(\mathbf{y}) \end{aligned}$$

and dividing by the constant $\beta - \alpha$ validates the claim $g(\mathbf{y}) \geq g(\mathbf{z})$. The claim $f(\mathbf{y}) \leq f(\mathbf{z})$ is proved by interchanging the roles of $f(\mathbf{x})$ and $g(\mathbf{x})$ and considering the functions $g(\mathbf{x}) + \alpha^{-1}f(\mathbf{x})$ and $g(\mathbf{x}) + \beta^{-1}f(\mathbf{x})$. ■

It is now fairly easy to prove a version of global convergence for the penalty method.

Proposition 13.2.2 *Suppose that both the objective function $f(\mathbf{x})$ and the penalty function $p(\mathbf{x})$ are continuous on \mathbb{R}^n and that the penalized functions $h_m(\mathbf{x}) = f(\mathbf{x}) + \lambda_m p(\mathbf{x})$ are coercive on \mathbb{R}^n . Then one can extract a corresponding sequence of minimum points \mathbf{x}_m such that $f(\mathbf{x}_m) \leq f(\mathbf{x}_{m+1})$. Furthermore, any cluster point of this sequence resides in the feasible region $C = \{\mathbf{x} : p(\mathbf{x}) = 0\}$ and attains the minimum value of $f(\mathbf{x})$ there. Finally, if $f(\mathbf{x})$ is coercive and possesses a unique minimum point in C , then the sequence \mathbf{x}_m converges to that point.*

Proof: In view of the coerciveness assumption, the minimum points \mathbf{x}_m exist. Proposition 13.2.1 confirms the ascent property. Now suppose that \mathbf{z} is a cluster point of the sequence \mathbf{x}_m and \mathbf{y} is any point in C . If we take limits in the inequalities

$$f(\mathbf{y}) = h_m(\mathbf{y}) \geq h_m(\mathbf{x}_m) \geq f(\mathbf{x}_m)$$

along the subsequence \mathbf{x}_{m_l} tending to \mathbf{z} , then the inequality $f(\mathbf{y}) \geq f(\mathbf{z})$ follows. Furthermore, because the λ_m tend to infinity, the bound

$$\limsup_{l \rightarrow \infty} \lambda_{m_l} p(\mathbf{x}_{m_l}) \leq f(\mathbf{y}) - \lim_{l \rightarrow \infty} f(\mathbf{x}_{m_l}) = f(\mathbf{y}) - f(\mathbf{z})$$

can only hold if $p(\mathbf{z}) = \lim_{l \rightarrow \infty} p(\mathbf{x}_{m_l}) = 0$.

If $f(\mathbf{x})$ possesses a unique minimum point \mathbf{y} in C , then to prove that \mathbf{x}_m converges to \mathbf{y} , it suffices to prove that \mathbf{x}_m is bounded. Given that $f(\mathbf{x})$ is coercive, it is possible to choose r so that $f(\mathbf{x}) > f(\mathbf{y})$ for all \mathbf{x} with $\|\mathbf{x}\| \geq r$. The assumption $\|\mathbf{x}_m\| \geq r$ consequently implies

$$h_m(\mathbf{x}_m) \geq f(\mathbf{x}_m) > f(\mathbf{y}) = h_m(\mathbf{y}),$$

which contradicts the assumption that \mathbf{x}_m minimizes $h_m(\mathbf{x})$. Hence, all \mathbf{x}_m satisfy $\|\mathbf{x}_m\| < r$. ■

Here is the corresponding result for the barrier method.

Proposition 13.2.3 *Suppose the real-valued function $f(\mathbf{x})$ is continuous on the bounded open set $U \subset \mathbb{R}^n$ and its closure V . Also suppose the barrier function $b(\mathbf{x})$ is continuous and coercive on U . If the tuning constants μ_m decrease to 0, then the linear combinations $h_m(\mathbf{x}) = f(\mathbf{x}) + \mu_m b(\mathbf{x})$ attain their minima at a sequence of points \mathbf{x}_m in U satisfying the descent property $f(\mathbf{x}_{m+1}) \leq f(\mathbf{x}_m)$. Furthermore, any cluster point of the sequence furnishes the minimum value of $f(\mathbf{x})$ on V . If the minimum point of $f(\mathbf{x})$ in V is unique, then the sequence \mathbf{x}_m converges to this point.*

Proof: Each of the continuous functions $h_m(\mathbf{x})$ is coercive on U , being the sum of a coercive function and a function bounded below. Therefore, the sequence \mathbf{x}_m exists. An appeal to Proposition 13.2.1 establishes the descent property. If \mathbf{z} is a cluster point of \mathbf{x}_m and \mathbf{x} is any point of U , then taking limits in the inequality

$$f(\mathbf{x}_m) + \mu_m b(\mathbf{x}_m) \leq f(\mathbf{x}) + \mu_m b(\mathbf{x})$$

along the relevant subsequence \mathbf{x}_{m_l} produces

$$f(\mathbf{z}) \leq \lim_{l \rightarrow \infty} f(\mathbf{x}_{m_l}) + \limsup_{l \rightarrow \infty} \mu_{m_l} b(\mathbf{x}_{m_l}) \leq f(\mathbf{x}).$$

It follows that $f(\mathbf{z}) \leq f(\mathbf{x})$ for every \mathbf{x} in V as well. If the minimum point of $f(\mathbf{x})$ on V is unique, then every cluster point of the bounded sequence \mathbf{x}_m coincides with this point. Hence, the sequence itself converges to the point. ■

Despite the elegance of the penalty and barrier methods, they suffer from three possible defects. First, they are predicated on finding the minimum point of the surrogate function for each value of the tuning constant. This entails iterations within iterations. Second, there is no obvious prescription for deciding how fast to send the tuning constants to their limits. Third, too large a value of λ_m in the penalty method or too small a value μ_m in the barrier method can lead to numerical instability in finding \mathbf{x}_m .

13.3 An Adaptive Barrier Method

The standard convex programming problem involves minimizing a convex function $f(\mathbf{x})$ subject to affine equality constraints $\mathbf{a}_i^* \mathbf{x} - b_i = 0$ for $1 \leq i \leq p$ and convex inequality constraints $h_j(\mathbf{x}) \leq 0$ for $1 \leq j \leq q$. This formulation renders the feasible region convex. To avoid distracting negative signs in this section, we will replace the constraint $h_j(\mathbf{x}) \leq 0$ by the constraint $v_j(\mathbf{x}) \geq 0$ for $v_j(\mathbf{x}) = -h_j(\mathbf{x})$. In the logarithmic barrier method, we define the barrier function

$$b(\mathbf{x}) = \sum_{j=1}^q \ln v_j(\mathbf{x}) \quad (13.2)$$

and optimize $g_m(\mathbf{x}) = f(\mathbf{x}) + \mu_m b(\mathbf{x})$ subject to the equality constraints. The presence of the barrier term $\ln v_j(\mathbf{x})$ keeps an initially inactive constraint $v_j(\mathbf{x})$ inactive throughout the search. Proposition 13.2.3 demonstrates convergence under specific hypotheses.

One way of improving the barrier method is to change the barrier constant as the iterations proceed [41, 162, 253]. This sounds vague, but matters simplify enormously if we view the construction of an adaptive barrier method from the perspective of the MM algorithm. Consider the following inequalities

$$\begin{aligned} & -v_j(\mathbf{x}_m) \ln v_j(\mathbf{x}) + v_j(\mathbf{x}_m) \ln v_j(\mathbf{x}_m) + dv_j(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) \\ \geq & -\frac{v_j(\mathbf{x}_m)}{v_j(\mathbf{x}_m)} [v_j(\mathbf{x}) - v_j(\mathbf{x}_m)] + dv_j(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) \\ = & -v_j(\mathbf{x}) + v_j(\mathbf{x}_m) + dv_j(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) \\ \geq & 0 \end{aligned} \quad (13.3)$$

based on the concavity of the functions $\ln y$ and $v_j(\mathbf{x})$. Because equality holds throughout when $\mathbf{x} = \mathbf{x}_m$, we have identified a novel function majorizing 0 and incorporating a barrier for $v_j(\mathbf{x})$. (Such functions are known as Bregman distances in the literature [24].) The significance of this discovery is that the surrogate function

$$\begin{aligned} g(\mathbf{x} \mid \mathbf{x}_m) &= f(\mathbf{x}) - \gamma \sum_{j=1}^q v_j(\mathbf{x}_m) \ln v_j(\mathbf{x}) \\ &+ \gamma \sum_{j=1}^q dv_j(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) \end{aligned} \quad (13.4)$$

majorizes $f(\mathbf{x})$ up to an irrelevant additive constant. Here γ is a fixed positive constant. Minimization of the surrogate function drives $f(\mathbf{x})$ downhill while keeping the inequality constraints inactive. In the limit, one or more of the inequality constraints may become active.

Because minimization of the surrogate function $g(\mathbf{x} \mid \mathbf{x}_m)$ cannot be accomplished in closed form, we must revert to the MM gradient algorithm. In performing one step of Newton's method, we need the first and second differentials

$$\begin{aligned} dg(\mathbf{x}_m \mid \mathbf{x}_m) &= df(\mathbf{x}_m) \\ d^2g(\mathbf{x}_m \mid \mathbf{x}_m) &= d^2f(\mathbf{x}_m) - \gamma \sum_{j=1}^q d^2v_j(\mathbf{x}_m) \\ &\quad + \gamma \sum_{j=1}^q \frac{1}{v_j(\mathbf{x}_m)} \nabla v_j(\mathbf{x}_m) dv_j(\mathbf{x}_m). \end{aligned}$$

In view of the convexity of $f(\mathbf{x})$ and the concavity of the $v_j(\mathbf{x})$, it is obvious that $d^2g(\mathbf{x}_m \mid \mathbf{x}_m)$ is positive semidefinite. It can be positive definite even if $d^2f(\mathbf{x}_m)$ is not.

As a safeguard in Newton's method, it is always a good idea to contract any proposed step so that simultaneously $f(\mathbf{x}_{m+1}) < f(\mathbf{x}_m)$ and $v_j(\mathbf{x}_{m+1}) \geq \delta v_j(\mathbf{x}_m)$ for all j and a small δ such as 0.1. It is also prudent to guard against ill conditioning of the matrix $d^2g(\mathbf{x}_m \mid \mathbf{x}_m)$ as a boundary $v_j(\mathbf{x}) = 0$ is approached and the multiplier $v_j(\mathbf{x}_m)^{-1}$ tends to ∞ . If one inverts $d^2g(\mathbf{x}_m \mid \mathbf{x}_m)$ or a bordered version of it by sweeping [170], then ill conditioning can be monitored as successive diagonal entries are swept. When the j th diagonal entry is dangerously close to 0, a small positive ϵ can be added to it just prior to sweeping. This apparently ad hoc remedy corresponds to adding the penalty $\frac{\epsilon}{2}(x_j - x_{mj})^2$ to $g(\mathbf{x} \mid \mathbf{x}_m)$. Although this action does not compromise the descent property, it does attenuate the parameter increment along the j th coordinate.

The surrogate function (13.4) does not exhaust the possibilities for majorizing the objective function. If we replace the concave function $\ln y$ by the concave function $-y^{-\alpha}$ in our derivation (13.3), then we can construct for each $\alpha > 0$ and β the alternative surrogate

$$\begin{aligned} g(\mathbf{x} \mid \mathbf{x}_m) &= f(\mathbf{x}) + \gamma \sum_{j=1}^q v_j(\mathbf{x}_m)^{\alpha+\beta} v_j(\mathbf{x})^{-\alpha} \\ &\quad + \gamma \alpha \sum_{j=1}^q v_j(\mathbf{x}_m)^{\beta-1} dv_j(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) \end{aligned} \quad (13.5)$$

majorizing $f(\mathbf{x})$ up to an irrelevant additive constant. This surrogate also exhibits an adaptive barrier that prevents the constraint $v_j(\mathbf{x})$ from becoming prematurely active. Imposing the condition $\alpha + \beta > 0$ is desirable because we want a barrier to relax as its boundary is approached. For this particular surrogate, straightforward differentiation yields

$$dg(\mathbf{x}_m \mid \mathbf{x}_m) = df(\mathbf{x}_m) \quad (13.6)$$

$$\begin{aligned}
 d^2g(\mathbf{x}_m | \mathbf{x}_m) &= d^2f(\mathbf{x}_m) - \gamma\alpha \sum_{j=1}^q v_j(\mathbf{x}_m)^{\beta-1} d^2v_j(\mathbf{x}_m) \quad (13.7) \\
 &\quad + \gamma\alpha(\alpha+1) \sum_{j=1}^q v_j(\mathbf{x}_m)^{\beta-2} \nabla v_j(\mathbf{x}_m) dv_j(\mathbf{x}_m).
 \end{aligned}$$

Example 13.3.1 *A Geometric Programming Example*

Consider the typical geometric programming problem of minimizing

$$f(\mathbf{x}) = \frac{1}{x_1x_2x_3} + x_2x_3$$

subject to

$$v(\mathbf{x}) = 4 - 2x_1x_3 - x_1x_2 \geq 0$$

and positive values for the x_i . Making the change of variables $x_i = e^{y_i}$ transforms the problem into a convex program. With the choice $\gamma = 1$, the MM gradient algorithm with the exponential parameterization and the log surrogate (13.4) produces the iterates displayed in the top half of Table 13.1. In this case Newton's method performs well, and none of the safeguards is needed. The MM gradient algorithm with the power surrogate (13.5) does somewhat better. The results shown in the bottom half of Table 13.1 reflect the choices $\gamma = 1$, $\alpha = 1/2$, and $\beta = 1$. ■

In the presence of linear constraints, both updates for the adaptive barrier method rely on the quadratic approximation of the surrogate function $g(\mathbf{x} | \mathbf{x}_m)$ using the calculated first and second differentials. This quadratic approximation is then minimized subject to the equality constraints as prescribed in Example 5.2.6.

Example 13.3.2 *Linear Programming*

Consider the standard linear programming problem of minimizing $\mathbf{c}^* \mathbf{x}$ subject to $\mathbf{A}\mathbf{x} = \mathbf{b}$ and $\mathbf{x} \geq \mathbf{0}$ [87]. At iteration $m+1$ of the adaptive barrier method with the power surrogate (13.5), we minimize the quadratic approximation

$$\mathbf{c}^* \mathbf{x}_m + \mathbf{c}^*(\mathbf{x} - \mathbf{x}_m) + \frac{1}{2} \gamma \alpha (\alpha + 1) \sum_{j=1}^n x_{mj}^{\beta-2} (x_j - x_{mj})^2$$

to the surrogate subject to $\mathbf{A}(\mathbf{x} - \mathbf{x}_m) = \mathbf{0}$. Note here the application of the two identities (13.6) and (13.7). According to Example 5.2.6 and equation (10.5), this minimization problem has solution

$$\mathbf{x}_{m+1} = \mathbf{x}_m - [\mathbf{D}_m^{-1} - \mathbf{D}_m^{-1} \mathbf{A}^* (\mathbf{A} \mathbf{D}_m^{-1} \mathbf{A}^*)^{-1} \mathbf{A} \mathbf{D}_m^{-1}] \mathbf{c},$$

where \mathbf{D}_m is a diagonal matrix with j th diagonal entry $\gamma\alpha(\alpha+1)x_{mj}^{\beta-2}$. It is convenient here to take $\gamma\alpha(\alpha+1) = 1$ and to step half along the search

TABLE 13.1. Solution of a geometric programming problem
 Iterates for the log surrogate

Iteration m	$f(\mathbf{x}_m)$	x_{m1}	x_{m2}	x_{m3}
1	2.0000	1.0000	1.0000	1.0000
2	1.7299	1.4386	0.9131	0.6951
3	1.6455	1.6562	0.9149	0.6038
4	1.5993	1.7591	0.9380	0.5685
5	1.5700	1.8256	0.9554	0.5478
10	1.5147	1.9614	0.9903	0.5098
15	1.5034	1.9910	0.9977	0.5023
20	1.5008	1.9979	0.9995	0.5005
25	1.5002	1.9995	0.9999	0.5001
30	1.5000	1.9999	1.0000	0.5000
35	1.5000	2.0000	1.0000	0.5000

Iterates for the power surrogate

1	2.0000	1.0000	1.0000	1.0000
2	1.6478	1.5732	1.0157	0.6065
3	1.5817	1.7916	0.9952	0.5340
4	1.5506	1.8713	1.0011	0.5164
5	1.5324	1.9163	1.0035	0.5090
10	1.5040	1.9894	1.0011	0.5008
15	1.5005	1.9986	1.0002	0.5001
20	1.5001	1.9998	1.0000	0.5000
25	1.5000	2.0000	1.0000	0.5000

direction $\mathbf{x}_{m+1} - \mathbf{x}_m$ whenever necessary. The case $\beta = 0$ bears a strong resemblance to Karmarkar’s celebrated method of linear programming. ■

We now show that the MM algorithms based on the surrogates (13.4) and (13.5) converge under fairly natural conditions. In the interests of generality, we will not require the objective function $f(\mathbf{x})$ to be convex. However, we will retain the assumptions of linear equality constraints $\mathbf{Ax} = \mathbf{b}$ and concave inequality constraints $v_j(\mathbf{x}) \geq 0$. To carry out our agenda, we assume that (a) $f(\mathbf{x})$ and the constraint functions $v_j(\mathbf{x})$ are continuously differentiable, (b) $f(\mathbf{x})$ is coercive, and (c) the second differential of $-\sum_{j=1}^q v_j(\mathbf{x})$ is positive definite on the affine subspace $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$. For simplicity, the objective function and the constraint functions are defined throughout \mathbb{R}^n . Either algorithm starts with a feasible point with all inequality constraints inactive.

For a subset $S \subset \{1, \dots, q\}$, let M_S be the active manifold defined by the equalities $\mathbf{Ax} = \mathbf{b}$ and $v_j(\mathbf{x}) = 0$ for $j \in S$ and the inequalities $v_j(\mathbf{x}) > 0$ for $j \notin S$. If M_S is empty, then we can safely ignore it in the sequel. Let $P_S(\mathbf{x})$ denote the projection matrix satisfying $dv_j(\mathbf{x})P_S(\mathbf{x}) = \mathbf{0}^*$ for every

$j \in S$ and defined by

$$P_S(\mathbf{x}) = \mathbf{I} - dV_S(\mathbf{x})^* [dV_S(\mathbf{x})dV_S(\mathbf{x})^*]^{-1} dV_S(\mathbf{x}), \quad (13.8)$$

where $dV_S(\mathbf{x})$ consists of the row vectors $dv_j(\mathbf{x})$ with $j \in S$ stacked one atop another. For the matrix inverse appearing in equation (13.8) to make sense, the matrix $dV_S(\mathbf{x})$ should have full row rank. The matrix $P_S(\mathbf{x})$ projects a row vector onto the subspace perpendicular to the differentials $dv_j(\mathbf{x})$ of the active constraints. For reasons that will become clear later, we insist that $\mathbf{A}P_S(\mathbf{x})$ have full row rank for each nonempty manifold M_S . When S is the empty set, we interpret $P_S(\mathbf{x})$ as the identity matrix \mathbf{I} .

We will call a point $\mathbf{x} \in M_S$ a stationary point if it satisfies the multiplier rule

$$df(\mathbf{x}) + \boldsymbol{\lambda}^* \mathbf{A} - \sum_{j \in S} \mu_j dv_j(\mathbf{x}) = \mathbf{0}^* \quad (13.9)$$

for some vector $\boldsymbol{\lambda}$ and collection of nonnegative coefficients μ_j . According to Proposition 6.5.3, a stationary point furnishes a global minimum of $f(\mathbf{x})$ when $f(\mathbf{x})$ is convex. We will assume that each manifold M_S possesses at most a finite number of stationary points. This is certainly the case when $f(\mathbf{x})$ is strictly convex, but it can also hold for linear or even non-convex objective functions [87].

Proposition 13.3.1 *Under the conditions just sketched, the adaptive barrier algorithm based on either the surrogate function (13.4) or the surrogate function (13.5) with $\beta = 1$ converges to a stationary point of $f(\mathbf{x})$. If $f(\mathbf{x})$ is convex, then the algorithms converge to the unique global minimum \mathbf{y} of $f(\mathbf{x})$ subject to the constraints.*

Proof: For the sake of brevity, we consider only the surrogate function (13.4). The coerciveness assumption guarantees that $f(\mathbf{x})$ possesses a minimum and that all iterates of a descent algorithm remain within a compact set. Because $g(\mathbf{x} \mid \mathbf{x}_m)$ majorizes $f(\mathbf{x})$, it is coercive and attains its minimum value as well. Unless $f(\mathbf{x})$ is convex, the minimum point \mathbf{x}_{m+1} of $g(\mathbf{x} \mid \mathbf{x}_m)$ may fail to be unique. When $f(\mathbf{x})$ is convex, assumption (c) implies that the quadratic form

$$\begin{aligned} \mathbf{u}^* d^2 g(\mathbf{x} \mid \mathbf{x}_m) \mathbf{u} &= \mathbf{u}^* d^2 f(\mathbf{x}) \mathbf{u} - \gamma \sum_{j=1}^q \frac{v_j(\mathbf{x}_m)}{v(\mathbf{x})} \mathbf{u}^* d^2 v_j(\mathbf{x}) \mathbf{u} \\ &\quad + \gamma \sum_{j=1}^q \frac{v_j(\mathbf{x}_m)}{v_j(\mathbf{x})^2} [dv_j(\mathbf{x}) \mathbf{u}]^2 \end{aligned}$$

is positive whenever $\mathbf{u} \neq \mathbf{0}$. Hence, $g(\mathbf{x} \mid \mathbf{x}_m)$ is strictly convex on the affine subspace $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$.

Given these preliminaries, our attack is based on taking limits in the stationarity equation

$$\begin{aligned} \mathbf{0}^* &= df(\mathbf{x}_{m+1}) + \boldsymbol{\lambda}_m^* A \\ &\quad - \gamma \sum_{j=1}^q \left[\frac{v_j(\mathbf{x}_m)}{v_j(\mathbf{x}_{m+1})} dv_j(\mathbf{x}_{m+1}) - dv_j(\mathbf{x}_m) \right] \end{aligned} \quad (13.10)$$

satisfied by $g(\mathbf{x} \mid \mathbf{x}_m)$ and recovering the Lagrange multiplier rule satisfied by $f(\mathbf{x})$. If we are to be successful in this regard, then we must show that

$$\lim_{m \rightarrow \infty} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| = 0. \quad (13.11)$$

Suppose the contrary is true. Then there exists a subsequence \mathbf{x}_{m_k} such that

$$\liminf_{k \rightarrow \infty} \|\mathbf{x}_{m_k+1} - \mathbf{x}_{m_k}\| > 0.$$

Invoking compactness and passing to a subsubsequence if necessary, we can also assume that $\lim_{k \rightarrow \infty} \mathbf{x}_{m_k} = \mathbf{u}$ and $\lim_{k \rightarrow \infty} \mathbf{x}_{m_k+1} = \mathbf{w}$ with $\mathbf{u} \neq \mathbf{w}$. In view of the inequalities (13.3) and $g(\mathbf{x}_{m+1} \mid \mathbf{x}_m) \leq g(\mathbf{x}_m \mid \mathbf{x}_m)$ and the concavity of $v_j(\mathbf{x})$ and $\ln t$, we deduce the further inequalities

$$\begin{aligned} 0 &\leq \gamma \sum_{j=1}^q [v_j(\mathbf{x}_{m_k}) - v_j(\mathbf{x}_{m_k+1}) + dv_j(\mathbf{x}_{m_k})(\mathbf{x}_{m_k+1} - \mathbf{x}_{m_k})] \\ &\leq \gamma \sum_{j=1}^q \left[v_j(\mathbf{x}_{m_k}) \ln \frac{v_j(\mathbf{x}_{m_k})}{v_j(\mathbf{x}_{m_k+1})} + dv_j(\mathbf{x}_{m_k})(\mathbf{x}_{m_k+1} - \mathbf{x}_{m_k}) \right] \\ &= g(\mathbf{x}_{m_k+1} \mid \mathbf{x}_{m_k}) - f(\mathbf{x}_{m_k+1}) - g(\mathbf{x}_{m_k} \mid \mathbf{x}_{m_k}) + f(\mathbf{x}_{m_k}) \\ &\leq f(\mathbf{x}_{m_k}) - f(\mathbf{x}_{m_k+1}). \end{aligned}$$

Given that $f(\mathbf{x}_m)$ is bounded and decreasing, in the limit the difference $f(\mathbf{x}_{m_k}) - f(\mathbf{x}_{m_k+1})$ tends to 0. It follows that

$$\gamma \sum_{j=1}^q [v_j(\mathbf{u}) - v_j(\mathbf{w}) + dv_j(\mathbf{u})(\mathbf{w} - \mathbf{u})] = 0,$$

contradicting the strict concavity of the sum $\sum_{j=1}^q v_j(\mathbf{x})$ on the affine subspace $\{\mathbf{x} : \mathbf{Ax} = \mathbf{b}\}$ and the hypothesis $\mathbf{u} \neq \mathbf{w}$.

Because the iterates \mathbf{x}_m all belong to the same compact set, the proposition can be proved by demonstrating that every convergent subsequence \mathbf{x}_{m_k} converges to the same stationary point \mathbf{y} . Consider such a subsequence with limit \mathbf{z} . Let us divide the constraint functions $v_j(\mathbf{x})$ into those that are active at \mathbf{z} and those that are inactive at \mathbf{z} . In the former case, we take $j \in S$, and in the latter case, we take $j \in S^c$. In a moment we will

demonstrate that \mathbf{z} is a stationary point of M_S . Because by hypothesis there are only a finite number of manifolds M_S and only a finite number of stationary points per manifold, Proposition 12.4.1 implies that all cluster points coincide and that the full sequence \mathbf{x}_m tends to a limit \mathbf{y} . To finish the proof, we must demonstrate that \mathbf{y} satisfies the multiplier rule for a constrained minimum. This last step is accomplished by taking limits in equality (13.10), assuming that λ_m and the ratios $v_j(\mathbf{x}_m)/v_j(\mathbf{x}_{m+1})$ all have limits. To avoid breaking the flow of our argument, we defer proof of these assertions. If $v_j(\mathbf{y}) > 0$, then it is obvious that $v_j(\mathbf{x}_m)/v_j(\mathbf{x}_{m+1})$ tends to 1, corresponding to a multiplier $\mu_j = 0$ for the inactive constraint j . If $v_j(\mathbf{y}) = 0$, then we must show that the limit of $v_j(\mathbf{x}_m)/v_j(\mathbf{x}_{m+1})$ exceeds 1. Otherwise, the multiplier μ_j is negative. But this limit relationship is valid because $v_j(\mathbf{x}_{m+1}) \leq v_j(\mathbf{x}_m)$ must hold for infinitely many m in order for $v_j(\mathbf{x}_m)$ to tend to 0.

We now return to the question of whether the limit \mathbf{z} of the convergent subsequence \mathbf{x}_{m_k} is a stationary point. To demonstrate that the subsequence λ_{m_k} converges, we multiply equation (13.10) on the right by the matrix $P_S(\mathbf{x}_{m_k+1})\mathbf{A}^*$ and solve for λ_{m_k} . This is possible because the matrix

$$B(\mathbf{x}_{m_k}) = \mathbf{A}P_S(\mathbf{x}_{m_k+1})\mathbf{A}^* = \mathbf{A}P_S(\mathbf{x}_{m_k+1})P_S(\mathbf{x}_{m_k+1})^*\mathbf{A}^*$$

has full rank by assumption. Since $P_S(\mathbf{x})$ annihilates $dv_j(\mathbf{x})$ for $j \in S$, and since \mathbf{x}_{m_k+1} converges to \mathbf{z} , a brief calculation shows that

$$\lambda^* = \lim_{k \rightarrow \infty} \lambda_{m_k} = -df(\mathbf{z})P_S(\mathbf{z})\mathbf{A}^*B(\mathbf{z})^{-1}.$$

To prove that the ratio $v_j(\mathbf{x}_{m_k})/v_j(\mathbf{x}_{m_k+1})$ has a limit for $j \in S$, we multiply equation (13.10) on the right by the matrix-vector product

$$P_{S-j}(\mathbf{x}_{m_k+1})\nabla v_j(\mathbf{x}_{m_k+1}),$$

where $S-j = S \setminus \{j\}$. This action annihilates all $dv_i(\mathbf{x}_{m_k+1})$ with $i \in S-j$ and makes it possible to express

$$\lim_{k \rightarrow \infty} \frac{v_j(\mathbf{x}_{m_k})}{v_j(\mathbf{x}_{m_k+1})} = \frac{[df(\mathbf{z}) + \lambda^*A + \gamma dv_j(\mathbf{z})]P_{S-j}(\mathbf{z})\nabla v_j(\mathbf{z})}{\gamma dv_j(\mathbf{z})P_{S-j}(\mathbf{z})\nabla v_j(\mathbf{z})}.$$

Note that the denominator $\gamma dv_j(\mathbf{z})P_{S-j}(\mathbf{z})\nabla v_j(\mathbf{z}) > 0$ because $\gamma > 0$ and $dV_{S-j}(\mathbf{z})$ has full row rank. Given these results, we can legitimately take limits in equation (13.10) along the given subsequence and recover the multiplier rule (13.9) at \mathbf{z} .

Now that we have demonstrated that \mathbf{x}_m tends to a unique limit \mathbf{y} , we can show that λ_m and the ratios $v_j(\mathbf{x}_m)/v_j(\mathbf{x}_{m+1})$ tend to well-defined limits by the logic employed with the subsequence \mathbf{x}_{m_k} . As noted earlier, this permits us to take limits in equation (13.10) and recover the multiplier

rule at \mathbf{y} . When $f(\mathbf{x})$ is convex, \mathbf{y} furnishes the global minimum. If there is another global minimum \mathbf{w} , then the entire line segment between \mathbf{w} and \mathbf{y} consists of minimum points. This contradicts the assumption that there are at most a finite number of stationary points throughout the feasible region. Hence, the minimum point \mathbf{y} is unique. ■

The next example illustrates that the local rate of convergence can be linear even when one of the constraints $v_i(\mathbf{x}) \geq 0$ is active at the minimum.

Example 13.3.3 *Convergence for the Multinomial Distribution*

As pointed out in Example 1.4.2, the loglikelihood for a multinomial distribution with d categories reduces to $\sum_{i=1}^d n_i \ln p_i$, where n_i is the observed number of counts in category i and p_i is the probability attached to category i . Maximizing the loglikelihood subject to the constraints $p_i \geq 0$ and $\sum_{i=1}^d p_i = 1$ gives the explicit maximum likelihood estimates $p_i = n_i/n$ for n trials. To compute the maximum likelihood estimates iteratively using the surrogate function (13.4), we find a stationary point of the Lagrangian

$$-\sum_{i=1}^d n_i \ln p_i - \gamma \sum_{i=1}^d p_{mi} \ln p_i + \gamma \sum_{i=1}^d (p_i - p_{mi}) + \lambda \left(\sum_{i=1}^d p_i - 1 \right).$$

Setting the i th partial derivative of the Lagrangian equal to 0 gives

$$-\frac{n_i}{p_i} - \frac{\gamma p_{mi}}{p_i} + \gamma + \lambda = 0. \tag{13.12}$$

Multiplying equation (13.12) by p_i , summing on i , and solving for λ yield $\lambda = n$. Substituting this value back in equation (13.12) produces

$$p_{m+1,i} = \frac{n_i + \gamma p_{mi}}{n + \gamma}.$$

At first glance it is not obvious that p_{mi} tends to n_i/n , but the algebraic rearrangement

$$\begin{aligned} p_{m+1,i} - \frac{n_i}{n} &= \frac{n_i + \gamma p_{mi}}{n + \gamma} - \frac{n_i}{n} \\ &= \frac{\gamma}{n + \gamma} \left(p_{mi} - \frac{n_i}{n} \right) \end{aligned}$$

shows that p_{mi} approaches n_i/n at the linear rate $\gamma/(n + \gamma)$. This is true regardless of whether $n_i/n = 0$ or $n_i/n > 0$. ■

13.4 Imposition of a Prior in EM Clustering

Priors imposed in Bayesian models play out as penalties in maximum a posteriori estimation. As an example, consider the EM clustering model studied in Sect. 9.5. There we postulated that each normally distributed

cluster had the same variance matrix. Relaxing this assumption sometimes causes the likelihood to become unbounded. Imposing a prior improves inference and stabilizes numerical estimation of parameters [44]. Let us review the derivation of the EM algorithm with these benefits in mind. The form of the EM updates

$$\pi_{n+1,j} = \frac{1}{m} \sum_{i=1}^m w_{ij}$$

for the admixture proportions π_j depend only on Bayes' rule and is valid regardless of the particular cluster densities. Here w_{ij} is the posterior probability that observation i comes from cluster j . To estimate the cluster means and common variance, we formed the surrogate function

$$\begin{aligned} & Q(\{\boldsymbol{\mu}_j, \boldsymbol{\Omega}_j\}_{j=1}^k \mid \{\boldsymbol{\mu}_{nj}, \boldsymbol{\Omega}_{nj}\}_{j=1}^k) \\ = & -\frac{1}{2} \sum_{j=1}^k \left(\sum_{i=1}^m w_{ij} \right) \ln \det \boldsymbol{\Omega}_j \\ & -\frac{1}{2} \sum_{j=1}^k \operatorname{tr} \left[\boldsymbol{\Omega}_j^{-1} \sum_{i=1}^m w_{ij} (\mathbf{y}_i - \boldsymbol{\mu}_j)(\mathbf{y}_i - \boldsymbol{\mu}_j)^* \right] \end{aligned}$$

with all $\boldsymbol{\Omega}_j = \boldsymbol{\Omega}$.

It is mathematically convenient to relax the common variance assumption and impose independent inverse Wishart priors on the different variance matrices $\boldsymbol{\Omega}_j$. In view of Problem 37 of Chap. 4, this amounts to adding the logprior

$$-\sum_{j=1}^k \left[\frac{a}{2} \ln \det \boldsymbol{\Omega}_j + \frac{b}{2} \operatorname{tr}(\boldsymbol{\Omega}_j^{-1} \mathbf{S}_j) \right]$$

to the surrogate function. Here the positive constants a and b and the positive definite matrices \mathbf{S}_j must be determined. Regardless of how these choices are made, we derive the usual EM updates

$$\boldsymbol{\mu}_{n+1,j} = \frac{1}{\sum_{i=1}^m w_{ij}} \sum_{i=1}^m w_{ij} \mathbf{y}_i$$

of the cluster means. Note that the constants a and b and the matrices \mathbf{S}_j have no influence on the weights w_{ij} computed via Bayes' rule.

The most natural choice is to take all \mathbf{S}_j equal to the sample variance matrix

$$\mathbf{S} = \frac{1}{m} \sum_{i=1}^m (\mathbf{y}_i - \bar{\mathbf{y}})(\mathbf{y}_i - \bar{\mathbf{y}})^*.$$

This choice is probably too diffuse, but it is better to err on the side of vagueness and avoid getting trapped at a local mode of the likelihood surface. In the absence of a prior, Example 4.7.6 implies that the EM update of Ω_j is

$$\tilde{\Omega}_{n+1,j} = \frac{1}{\sum_{i=1}^m w_{ij}} \sum_{i=1}^m w_{ij} (\mathbf{y}_i - \boldsymbol{\mu}_{n+1,j})(\mathbf{y}_i - \boldsymbol{\mu}_{n+1,j})^*$$

By the same reasoning, the maximum of the penalized surrogate function with respect to Ω_j is

$$\Omega_{n+1,j} = \frac{a}{a + \sum_{i=1}^m w_{ij}} \left(\frac{b}{a} \mathbf{S} \right) + \frac{\sum_{i=1}^m w_{ij}}{a + \sum_{i=1}^m w_{ij}} \tilde{\Omega}_{n+1,j}.$$

In other words, the penalized EM update is a convex combination of the standard EM update and the mode $\frac{b}{a} \mathbf{S}$ of the prior. Chen and Tan [44] tentatively recommend the choice $a = b = 2/\sqrt{m}$. As m tends to ∞ , the influence of the prior diminishes.

13.5 Model Selection and the Lasso

We now turn to penalized regression and continuous model selection. Our focus will be on the lasso penalty and its application in regression problems where the number of predictors p exceeds the number of cases n [45, 49, 229, 252, 256]. The lasso also finds applications in generalized linear models. In each of these contexts, let y_i be the response for case i , x_{ij} be the value of predictor j for case i , and β_j be the regression coefficient corresponding to predictor j . In practice one should standardize each predictor to have mean 0 and variance 1. Standardization puts all regression coefficients on a common scale as implicitly demanded by the lasso penalty.

The intercept α is ignored in the lasso penalty, whose strength is determined by the positive tuning constant λ . If $\boldsymbol{\theta} = (\alpha, \beta_1, \dots, \beta_p)^*$ is the parameter vector and $g(\boldsymbol{\theta})$ is the loss function ignoring the penalty, then the lasso minimizes the criterion

$$f(\boldsymbol{\theta}) = g(\boldsymbol{\theta}) + \lambda \sum_{j=1}^p |\beta_j|,$$

where $g(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (y_i - \mathbf{x}_i^* \boldsymbol{\theta})^2$ in ℓ_2 regression and $g(\boldsymbol{\theta}) = \sum_{i=1}^n |y_i - \mathbf{x}_i^* \boldsymbol{\theta}|$ in ℓ_1 regression. The penalty $\lambda \sum_j |\beta_j|$ shrinks each β_j toward the origin and tends to discourage models with large numbers of irrelevant predictors. The lasso penalty is more effective in this regard than the ridge penalty $\lambda \sum_j \beta_j^2$ because $|b|$ is much bigger than b^2 for small b .

Lasso penalized estimation raises two issues. First, what is the most effective method of minimizing the objective function $f(\boldsymbol{\theta})$? In the current

section we highlight the method of coordinate descent [55, 98, 100, 276]. Second, how does one choose the tuning constant λ ? The standard answer is cross-validation. Although this is a good reply, it does not resolve the problem of how to minimize average cross-validation error as measured by the loss function. Recall that in k -fold cross-validation, one divides the data into k equal batches (subsamples) and estimates parameters k times, leaving one batch out per time. The testing error (total loss) for each omitted batch is computed using the estimates derived from the remaining batches, and the cross-validation error $c(\lambda)$ is computed by averaging testing error across the k batches.

Unless carefully planned, evaluation of $c(\lambda)$ on a grid of points may be computationally costly, particularly if grid points occur near $\lambda = 0$. Because coordinate descent is fastest when λ is large and the vast majority of β_j are estimated as 0, it makes sense to start with a very large value and work downward. One advantage of this tactic is that parameter estimates for a given λ can be used as parameter starting values for the next lower λ . For the initial value of λ , the starting value $\boldsymbol{\theta} = \mathbf{0}$ is recommended. It is also helpful to set an upper bound on the number of active parameters allowed and abort downward sampling of λ when this bound is exceeded. Once a fine enough grid is available, visual inspection usually suggests a small interval flanking the minimum. Application of golden section search over the flanking interval will then quickly lead to the minimum.

Coordinate descent comes in several varieties. The standard version cycles through the parameters and updates each in turn. An alternative version is greedy and updates the parameter giving the largest decrease in the objective function. Because it is impossible to tell in advance the extent of each decrease, the greedy version uses the surrogate criterion of steepest descent. In other words, for each parameter we compute forward and backward directional derivatives and update the parameter with the most negative directional derivative, either forward or backward. The overhead of keeping track of the directional derivative works to the detriment of the greedy method. For ℓ_1 regression, the overhead is relatively light, and greedy coordinate descent converges faster than cyclic coordinate descent.

Although the lasso penalty is nondifferentiable, it does possess directional derivatives along each forward or backward coordinate direction. For instance, if \mathbf{e}_j is the coordinate direction along which β_j varies, then

$$d_{\mathbf{e}_j} f(\boldsymbol{\theta}) = \lim_{t \downarrow 0} \frac{f(\boldsymbol{\theta} + t\mathbf{e}_j) - f(\boldsymbol{\theta})}{t} = d_{\mathbf{e}_j} g(\boldsymbol{\theta}) + \begin{cases} \lambda & \beta_j \geq 0 \\ -\lambda & \beta_j < 0, \end{cases}$$

and

$$d_{-\mathbf{e}_j} f(\boldsymbol{\theta}) = \lim_{t \downarrow 0} \frac{f(\boldsymbol{\theta} - t\mathbf{e}_j) - f(\boldsymbol{\theta})}{t} = d_{-\mathbf{e}_j} g(\boldsymbol{\theta}) + \begin{cases} -\lambda & \beta_j > 0 \\ \lambda & \beta_j \leq 0. \end{cases}$$

In ℓ_1 regression, the loss function is also nondifferentiable, and a brief calculation shows that the coordinate directional derivatives are

$$d_{\mathbf{e}_j} \sum_{i=1}^n |y_i - \mathbf{x}_i^* \boldsymbol{\theta}| = \sum_{i=1}^n \begin{cases} -x_{ij} & y_i - \mathbf{x}_i^* \boldsymbol{\theta} > 0 \\ x_{ij} & y_i - \mathbf{x}_i^* \boldsymbol{\theta} < 0 \\ |x_{ij}| & y_i - \mathbf{x}_i^* \boldsymbol{\theta} = 0 \end{cases}$$

and

$$d_{-\mathbf{e}_j} \sum_{i=1}^n |y_i - \mathbf{x}_i^* \boldsymbol{\theta}| = \sum_{i=1}^n \begin{cases} x_{ij} & y_i - \mathbf{x}_i^* \boldsymbol{\theta} > 0 \\ -x_{ij} & y_i - \mathbf{x}_i^* \boldsymbol{\theta} < 0 \\ |x_{ij}| & y_i - \mathbf{x}_i^* \boldsymbol{\theta} = 0 \end{cases}$$

with predictor vector $\mathbf{x}_i^* = (1, \mathbf{z}_i^*)$ for case i . Fortunately, when a function is differentiable, its directional derivative along \mathbf{e}_j coincides with its ordinary partial derivative, and its directional derivative along $-\mathbf{e}_j$ coincides with the negative of its ordinary partial derivative.

When we visit parameter β_j in cyclic coordinate descent, we evaluate $d_{\mathbf{e}_j} f(\boldsymbol{\theta})$ and $d_{-\mathbf{e}_j} f(\boldsymbol{\theta})$. If both are nonnegative, then we skip the update for β_j . This decision is defensible when $g(\boldsymbol{\theta})$ is convex because the sign of a directional derivative fully determines whether improvement can be made in that direction. If either directional derivative is negative, then we must solve for the minimum in that direction. When the current slope parameter β_j is parked at 0 and the partial derivative $\frac{\partial}{\partial \beta_j} g(\boldsymbol{\theta})$ exists,

$$d_{\mathbf{e}_j} f(\boldsymbol{\theta}) = \frac{\partial}{\partial \beta_j} g(\boldsymbol{\theta}) + \lambda, \quad d_{-\mathbf{e}_j} f(\boldsymbol{\theta}) = -\frac{\partial}{\partial \beta_j} g(\boldsymbol{\theta}) + \lambda.$$

Hence, β_j moves to the right if $\frac{\partial}{\partial \beta_j} g(\boldsymbol{\theta}) < -\lambda$, to the left if $\frac{\partial}{\partial \beta_j} g(\boldsymbol{\theta}) > \lambda$, and stays fixed otherwise. In underdetermined problems with just a few relevant predictors, most updates are skipped, and the parameters never budge from their starting values of 0. This simple fact plus the complete absence of matrix operations explains the speed of coordinate descent. It inherits its numerical stability from the descent property of each update.

13.6 Lasso Penalized ℓ_1 Regression

In lasso constrained ℓ_1 regression, greedy coordinate descent is quick because directional derivatives are trivial to update. Indeed, if updating β_j does not alter the sign of the residual $y_i - \mathbf{x}_i^* \boldsymbol{\theta}$ for case i , then the contributions of case i to the various directional derivatives do not change. When the residual $y_i - \mathbf{x}_i^* \boldsymbol{\theta}$ changes sign, these contributions change by $\pm 2x_{ij}$. When a residual changes from 0 to nonzero or vice versa, the increment depends on the sign of the nonzero residual and the sign of x_{ij} .

Updating the value of the chosen parameter can be achieved by the nearly forgotten algorithm of Edgeworth [80, 81], which for a long time was

considered a competitor of least squares. Portnoy and Koenker [214] trace the history of the algorithm from Boscovich to Laplace to Edgeworth. It is fair to say that the algorithm has managed to cling to life despite decades of obscurity both before and after its rediscovery by Edgeworth.

To illustrate Edgeworth's algorithm in operation, consider minimizing the two-parameter model

$$g(\boldsymbol{\theta}) = \sum_{i=1}^n |y_i - \alpha - z_i \beta|$$

with a single slope β . To update α , we recall the well-known connection between ℓ_1 regression and medians and replace α for fixed β by the sample median of the numbers $v_i = y_i - z_i \beta$. This action drives $g(\boldsymbol{\theta})$ downhill. Updating β for α fixed depends on writing

$$g(\boldsymbol{\theta}) = \sum_{i=1}^n |z_i| \left| \frac{y_i - \alpha}{z_i} - \beta \right|,$$

sorting the numbers $v_i = (y_i - \alpha)/z_i$, and finding the weighted median with weight $w_i = |z_i|$ assigned to v_i . We replace β by the order statistic $v_{[i]}$ with weight $w_{[i]}$ whose index i satisfies

$$\sum_{j=1}^{i-1} w_{[j]} < \frac{1}{2} \sum_{j=1}^n w_{[j]}, \quad \sum_{j=1}^i w_{[j]} \geq \frac{1}{2} \sum_{j=1}^n w_{[j]}.$$

Problem 6 demonstrates that this choice is valid. Edgeworth's algorithm easily generalizes to multiple linear regression. Implementing the algorithm with a lasso penalty requires viewing the penalty terms as the absolute values of pseudo-residuals. Thus, we write

$$\lambda |\beta_j| = |y - \mathbf{x}^* \boldsymbol{\theta}|$$

by taking $y = 0$ and $x_k = \lambda 1_{\{k=j\}}$.

Two criticisms have been leveled at Edgeworth's algorithm. First, although it drives the objective function steadily downhill, it sometimes stalls at an inferior point. See Problem 8 for an example. The second criticism is that convergence often occurs in a slow seesaw pattern. These defects are not completely fatal. As late as 1978, Armstrong and Kung published a computer implementation of Edgeworth's algorithm in the journal *Applied Statistics* [4].

13.7 Lasso Penalized ℓ_2 Regression

In ℓ_2 regression with a lasso penalty, we minimize the objective function

$$f(\boldsymbol{\theta}) = \frac{1}{2} \sum_{i=1}^n (y_i - \alpha - \mathbf{z}_i^* \boldsymbol{\beta})^2 + \lambda \sum_{j=1}^p |\beta_j| = g(\boldsymbol{\theta}) + \lambda \sum_{j=1}^p |\beta_j|.$$

The update of the intercept parameter can be written as

$$\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n (y_i - \mathbf{z}_i^* \boldsymbol{\beta}) = \alpha - \frac{1}{n} \frac{\partial}{\partial \alpha} g(\boldsymbol{\theta}).$$

For the parameter β_k , there are separate solutions to the left and right of 0. These boil down to

$$\begin{aligned} \hat{\beta}_{k,-} &= \min \left\{ 0, \beta_k - \frac{\frac{\partial}{\partial \beta_k} g(\boldsymbol{\theta}) - \lambda}{\sum_{i=1}^n z_{ik}^2} \right\} \\ \hat{\beta}_{k,+} &= \max \left\{ 0, \beta_k - \frac{\frac{\partial}{\partial \beta_k} g(\boldsymbol{\theta}) + \lambda}{\sum_{i=1}^n z_{ik}^2} \right\}. \end{aligned}$$

The reader can check that only one of these two solutions can be nonzero. The partial derivatives

$$\frac{\partial}{\partial \alpha} g(\boldsymbol{\theta}) = - \sum_{i=1}^n r_i, \quad \frac{\partial}{\partial \beta_k} g(\boldsymbol{\theta}) = - \sum_{i=1}^n r_i z_{ik}$$

of $g(\boldsymbol{\theta})$ are easy to compute provided we keep track of all of the residuals $r_i = y_i - \alpha - \mathbf{z}_i^* \boldsymbol{\beta}$. The residual r_i starts with the value y_i and is reset to $r_i + \alpha - \hat{\alpha}$ when α is updated and to $r_i + z_{ij}(\beta_j - \hat{\beta}_j)$ when β_j is updated. Organizing all updates around residuals promotes fast evaluation of $g(\boldsymbol{\theta})$. At the expense of somewhat more complex code [99], a better tactic is to exploit the identity

$$\sum_{i=1}^n r_i z_{ik} = \sum_{i=1}^n y_i z_{ik} - \alpha \sum_{i=1}^n z_{ik} - \sum_{j:|\beta_j|>0} \left(\sum_{i=1}^n z_{ij} z_{ik} \right) \beta_j.$$

This representation suggests storing and reusing the inner products

$$\sum_{i=1}^n y_i z_{ik}, \quad \sum_{i=1}^n z_{ik}, \quad \sum_{i=1}^n z_{ij} z_{ik}$$

for the active predictors.

Example 13.7.1 Obesity and Gene Expression in Mice

Consider a genetics example involving gene expression levels and obesity in mice. Wang et al. [268] measured abdominal fat mass on $n = 311$ F2 mice (155 males and 156 females). The F2 mice were created by mating two inbred strains and then mating brother-sister pairs from the resulting offspring. Wang et al. [268] also recorded the expression levels in liver of $p = 23,388$ genes in each mouse. A reasonable model postulates

$$y_i = 1_{\{i \text{ male}\}} \alpha_1 + 1_{\{i \text{ female}\}} \alpha_2 + \sum_{j=1}^p x_{ij} \beta_j + \epsilon_i,$$

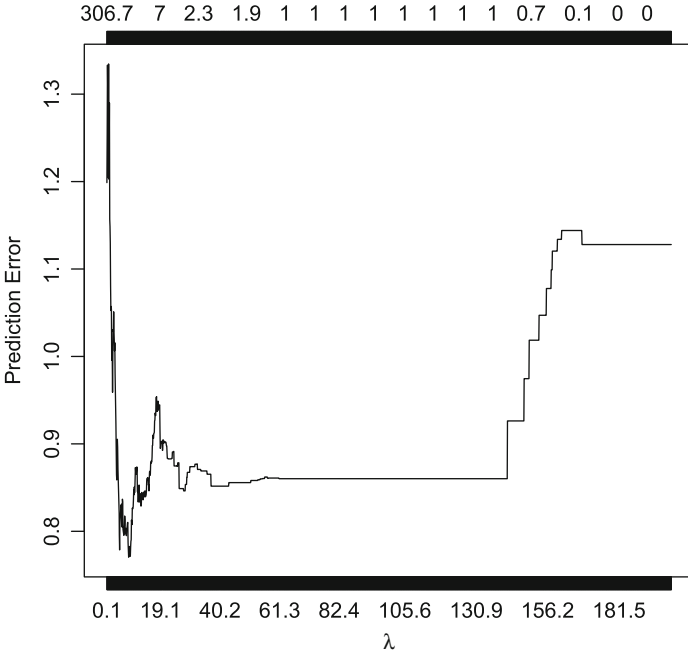


FIGURE 13.1. The cross-validation curve $c(\lambda)$ for obesity in mice

where y_i measures fat mass on mouse i , x_{ij} is the expression level of gene j in mouse i , and ϵ_i is random error. Since male and female mice exhibit across the board differences in size and physiology, it is prudent to estimate a different intercept for each sex. Figure 13.1 plots average prediction error as a function of λ (lower horizontal axis) and the average number of nonzero predictors (upper horizontal axis). Here we use ℓ_2 penalized regression and 10-fold cross-validation. Examination of the cross-validation curve $c(\lambda)$ over a fairly dense grid shows an optimal λ of 7.8 with 41 nonzero predictors. For ℓ_1 penalized regression, the optimal λ is around 3.5 with 77 nonzero predictors. The preferred ℓ_1 and ℓ_2 models share 27 predictors in common. Several of the genes identified are known or suspected to be involved in lipid metabolism, adipose deposition, and impaired insulin sensitivity in mice. More details can be found in the paper [276].

The tactics described for ℓ_2 regression carry over to generalized linear models. In this setting, the loss function $g(\theta)$ is the negative loglikelihood. In many cases, $g(\theta)$ is convex, and it is possible to determine whether progress can be made along a forward or backward coordinate direction without actually minimizing the objective function. It is clearly computationally beneficial to organize parameter updates by tracking the linear predictor $\alpha + z_i^* \beta$ of each case. Although we no longer have explicit solutions

to fall back on, the scoring algorithm serves as a substitute. Since it usually converges in a few iterations, the computational overhead of cyclic coordinate descent remains manageable.

13.8 Penalized Discriminant Analysis

Discriminant analysis is another attractive candidate for penalized estimation. In discriminant analysis with two categories, each case i is characterized by a feature vector \mathbf{z}_i and a category membership indicator y_i taking the values -1 or 1 . In the machine learning approach to discriminant analysis [231, 264], the hinge loss function $[1 - y_i(\alpha + \mathbf{z}_i^* \boldsymbol{\beta})]_+$ plays a prominent role. Here u_+ is shorthand for the convex function $\max\{u, 0\}$. Just as in ordinary regression, we can penalize the overall loss

$$g(\boldsymbol{\theta}) = \sum_{i=1}^n [1 - y_i(\alpha + \mathbf{z}_i^* \boldsymbol{\beta})]_+$$

by imposing a lasso or ridge penalty. Note that the linear regression function $h_i(\boldsymbol{\theta}) = \alpha + \mathbf{z}_i^* \boldsymbol{\beta}$ predicts either -1 or 1 . If $y_i = 1$ and $h_i(\boldsymbol{\theta})$ over-predicts in the sense that $h_i(\boldsymbol{\theta}) > 1$, then there is no loss. Similarly, if $y_i = -1$ and $h_i(\boldsymbol{\theta})$ under-predicts in the sense that $h_i(\boldsymbol{\theta}) < -1$, then there is no loss.

Most strategies for estimating $\boldsymbol{\theta}$ pass to the dual of the original minimization problem. A simpler strategy is to majorize each contribution to the loss by a quadratic and minimize the surrogate loss plus penalty [114]. A little calculus shows that $(u)_+$ is majorized at $u_m \neq 0$ by the quadratic

$$q(u \mid u_m) = \frac{1}{4|u_m|} (u + |u_m|)^2. \quad (13.13)$$

(See Problem 13.) In fact, this is the best quadratic majorizer of u_+ [62]. Both of the majorizations (8.12) and (13.13) have singularities at the point $u_m = 0$. One simple fix is to replace $|u_m|$ by $|u_m| + \epsilon$ wherever $|u_m|$ appears in a denominator in either formula. We recommend double precision arithmetic with $0 < \epsilon \leq 10^{-5}$. Problem 14 explores a more sophisticated remedy that replaces the functions $|u|$ and u_+ by differentiable approximations.

In any case, if we impose a ridge penalty, then the hinge majorization leads to a pure MM algorithm exploiting weighted least squares. Coordinate descent algorithms with a lasso or ridge penalty are also enabled by majorization, but each coordinate update merely decreases the objective function along the given coordinate direction. Fortunately, this drawback is outweighed by the gain in numerical simplicity in majorizing hinge loss. The decisions to use a lasso or ridge penalty and apply pure MM or coordinate descent with majorization will be dictated in practical problems by the number of potential predictors. If a lasso penalty is imposed to

eliminate irrelevant predictors, then cyclic coordinate descent is preferable, with the surrogate function substituting for the objective function in each parameter update.

In discriminant analysis with more than two categories, it is convenient to pass to ϵ -insensitive loss and multiple linear regression. The story is too long to tell here, but it is worth mentioning that the conjunction of a parsimonious loss function and efficient MM or coordinate descent algorithms produce some of the most effective discriminant analysis methods tested [169, 277].

13.9 Problems

1. In Example 13.2.1 prove directly that the solution displayed in equation (13.1) converges to the minimum point of $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2$ subject to the linear constraints $\mathbf{V}\boldsymbol{\beta} = \mathbf{d}$. (Hints: Assume that the matrix \mathbf{V} has full column rank and consult Example 5.2.6, Proposition 5.2.2, and Problem 10 of Chap. 11.)
2. Prove that the surrogate function (13.5) majorizes $f(\mathbf{x})$ up to an irrelevant additive constant.
3. The power plant production problem [226] involves minimizing

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(x_i), \quad f_i(x_i) = a_i x_i + \frac{1}{2} b_i x_i^2$$

subject to the constraints $0 \leq x_i \leq u_i$ for each i and $\sum_{i=1}^n x_i \geq d$. For plant i , x_i is the power output, u_i is the capacity, and $f_i(x_i)$ is the cost. The total demand is d , and the cost constants a_i and b_i are positive. This problem can be solved by the adaptive barrier algorithm. Program this algorithm and test it on a simple example with at least two power plants. Argue that the minimum is unique. Example 15.6.1 sketches another approach.

4. In Problem 3 investigate the performance of cyclic coordinate descent. Explain why it fails.
5. Implement and test the EM clustering algorithm with a Bayesian prior. Apply the algorithm to Fisher's classic iris data set. Fisher's data can be downloaded from the web. See the book [191] for commentary and further references.
6. Show that $\hat{\mu}$ minimizes $f(\mu) = \sum_{i=1}^n w_i |x_i - \mu|$ if and only if

$$\sum_{x_i < \hat{\mu}} w_i \leq \frac{1}{2} \sum_{i=1}^n w_i \quad \text{and} \quad \sum_{x_i \leq \hat{\mu}} w_i \geq \frac{1}{2} \sum_{i=1}^n w_i.$$

Assume that the weights w_i are positive. (Hint: Apply Proposition 6.5.2.)

7. Consider the piecewise linear function

$$f(\mu) = c\mu + \sum_{i=1}^n w_i |x_i - \mu|,$$

where the positive weights satisfy $\sum_{i=1}^n w_i = 1$ and the points satisfy $x_1 < x_2 < \dots < x_n$. Show that $f(\mu)$ has no minimum when $|c| > 1$. What happens when $c = 1$ or $c = -1$? This leaves the case $|c| < 1$. Show that a minimum occurs when

$$\sum_{x_i > \mu} w_i - \sum_{x_i \leq \mu} w_i \leq c \quad \text{and} \quad \sum_{x_i \geq \mu} w_i - \sum_{x_i < \mu} w_i \geq c.$$

(Hints: A crude plot of $f(\mu)$ might help. What conditions on the right-hand and left-hand derivatives of $f(\mu)$ characterize a minimum?)

8. Show that Edgeworth's algorithm [178] for ℓ_1 regression converges to an inferior point for the data values $(0.3, -1.0)$, $(-0.4, -0.1)$, $(-2.0, -2.9)$, $(-0.9, -2.4)$, and $(-1.1, 2.2)$ for the pairs (x_i, y_i) and parameter starting values $(\alpha, \beta) = (3.5, -1.0)$.
9. Implement and test greedy coordinate descent for lasso penalized ℓ_1 regression or cyclic coordinate descent for lasso penalized ℓ_2 regression.
10. In lasso penalized regression, suppose the convex loss function $g(\boldsymbol{\theta})$ is differentiable. A stationary point $\boldsymbol{\theta}$ of coordinate descent satisfies the conditions $d_{\mathbf{e}_j} f(\boldsymbol{\theta}) \geq 0$ and $d_{-\mathbf{e}_j} f(\boldsymbol{\theta}) \geq 0$ for all j . Here the intercept α varies along the coordinate direction \mathbf{e}_0 . Calculate the general directional derivative

$$d_{\mathbf{v}} f(\boldsymbol{\theta}) = \sum_j \frac{\partial}{\partial \theta_j} g(\boldsymbol{\theta}) v_j + \lambda \sum_{j>0} \begin{cases} v_j & \theta_j > 0 \\ -v_j & \theta_j < 0 \\ |v_j| & \theta_j = 0 \end{cases}$$

and show that

$$d_{\mathbf{v}} f(\boldsymbol{\theta}) = \sum_{v_j > 0} d_{\mathbf{e}_j} f(\boldsymbol{\theta}) v_j + \sum_{v_j < 0} d_{-\mathbf{e}_j} f(\boldsymbol{\theta}) |v_j|.$$

Conclude that every directional derivative is nonnegative at a stationary point. In view of Proposition 6.5.2, stationary points therefore coincide with minimum points. This result does not hold for lasso penalized ℓ_1 regression.

11. Show that the function $\|\mathbf{x}\|_0 = \sum_{i=1}^n 1_{\{x_i \neq 0\}}$ satisfies the properties:

- (a) $\|\mathbf{x}\|_0$ is nonnegative and equal to 0 if and only if $\mathbf{x} = \mathbf{0}$,
- (b) $\|\mathbf{x}\|_0 = \|-\mathbf{x}\|_0$,
- (c) $\|\mathbf{x} + \mathbf{y}\|_0 \leq \|\mathbf{x}\|_0 + \|\mathbf{y}\|_0$,
- (d) The function $\mathbf{x} \mapsto \|\mathbf{x}\|_0$ is lower semicontinuous.

What norm property fails?

12. For the ℓ_0 “norm” $\|\mathbf{x}\|_0$ defined in the previous problem, demonstrate that

$$\|\mathbf{x}\|_0 = \lim_{\epsilon \downarrow 0} \sum_{i=1}^n \frac{\ln\left(1 + \frac{|x_i|}{\epsilon}\right)}{\ln\left(1 + \frac{1}{\epsilon}\right)}.$$

Note that the same limit applies if one substitutes x_i^2 for $|x_i|$. Now prove the majorization

$$\ln(\epsilon + y) \leq \ln(\epsilon + y_m) + \frac{1}{\epsilon + y_m} (y - y_m)$$

for nonnegative scalars y and y_m , and show how it can be employed to majorize approximations to $\|\mathbf{x}\|_0$ based on the choices $|x_i|$ and x_i^2 . See the references [39, 88, 272] for applications to sparse estimation and machine learning.

13. Show that the function $u_+ = \max\{u, 0\}$ is majorized by the quadratic function (13.13) at a point $u_m \neq 0$. Why does it suffice to prove that u_+ and $q(u \mid u_m)$ have the same value and same derivative at u_m and $-u_m$? Also check that u_+^2 is majorized by u^2 for $u_m \geq 0$ and by $(u - u_m)^2$ for $u_m < 0$. (Hint: Draw rough graphs of u_+ and $q(u \mid u_m)$.)

14. For a small $\epsilon > 0$, the functions $\sqrt{u^2 + \epsilon} - \sqrt{\epsilon}$ and $\sqrt{u_+^2 + \epsilon} - \sqrt{\epsilon}$ are excellent differentiable approximations to $|u|$ and u_+ , respectively. Derive the majorizations

$$\begin{aligned} \sqrt{u^2 + \epsilon} - \sqrt{\epsilon} &\leq \sqrt{u_m^2 + \epsilon} - \sqrt{\epsilon} + \frac{1}{2\sqrt{u_m^2 + \epsilon}}(u^2 - u_m^2) \\ \sqrt{u_+^2 + \epsilon} - \sqrt{\epsilon} &\leq \begin{cases} \frac{1}{2\sqrt{\epsilon}}u^2 & u_m = 0 \\ \frac{\sqrt{u_m^2 + \epsilon} + \sqrt{\epsilon}}{4(u_m^2 + \epsilon)} \left[u + \frac{u_m^2}{(\sqrt{u_m^2 + \epsilon} + \sqrt{\epsilon})^2} \right]^2 & u_m < 0 \\ \frac{\sqrt{u_m^2 + \epsilon} - \sqrt{\epsilon}}{(u_m - u_m)^2}(u - u_m)^2 & u_m > 0, \end{cases} \end{aligned}$$

where in the last case r_m is the largest real root of the cubic equation $u^3 + 2u_m u^2 + u_m^2 u + 4\epsilon u_m = 0$. (Hints: In majorizing the approximation to u_+ , in each case assume that $q(u | u_m) = c(u - d)^2$. Choose c and d to give one or two tangency points.)

15. Implement and test one of the discriminant analysis algorithms that depend on quadratic majorization of hinge loss.
16. Nonnegative matrix factorization was introduced by Lee and Seung [174, 175] as an analog of principal components and vector quantization with applications in data compression and clustering. In mathematical terms, one approximates a matrix \mathbf{U} with nonnegative entries u_{ij} by a product $\mathbf{V}\mathbf{W}$ of two low-rank matrices with nonnegative entries v_{ij} and w_{ij} . If the entries u_{ij} are integers, then they can be viewed as realizations of independent Poisson random variables with means $\sum_k v_{ik} w_{kj}$. In this setting the loglikelihood is

$$L(\mathbf{V}, \mathbf{W}) = \sum_i \sum_j \left[u_{ij} \ln \left(\sum_k v_{ik} w_{kj} \right) - \sum_k v_{ik} w_{kj} \right].$$

Maximization with respect to \mathbf{V} and \mathbf{W} should lead to a good factorization. Lee and Seung construct a block ascent algorithm that hinges on the minorization

$$\ln \left(\sum_k v_{ik} w_{kj} \right) \geq \sum_k \frac{a_{ikj}^n}{b_{ij}^n} \ln \left(\frac{b_{ij}^n}{a_{ikj}^n} v_{ik} w_{kj} \right),$$

where

$$a_{ikj}^n = v_{ik}^n w_{kj}^n, \quad b_{ij}^n = \sum_k v_{ik}^n w_{kj}^n,$$

and n indicates the current iteration. Prove this minorization and derive the Lee-Seung algorithm with alternating multiplicative updates

$$v_{ik}^{n+1} = v_{ik}^n \frac{\sum_j u_{ij} \frac{w_{kj}^n}{b_{ij}^n}}{\sum_j w_{kj}^n}$$

and

$$w_{kj}^{n+1} = w_{kj}^n \frac{\sum_i u_{ij} \frac{v_{ik}^n}{b_{ij}^n}}{\sum_i v_{ik}^n}.$$

17. Continuing Problem 16, consider minimizing the squared Frobenius norm

$$\|\mathbf{U} - \mathbf{V}\mathbf{W}\|_F^2 = \sum_i \sum_j \left(u_{ij} - \sum_k v_{ik} w_{kj} \right)^2.$$

Demonstrate the majorization

$$\left(u_{ij} - \sum_k v_{ik} w_{kj}\right)^2 \leq \sum_k \frac{a_{ikj}^n}{b_{ij}^n} \left(u_{ij} - \frac{b_{ij}^n}{a_{ikj}^n} v_{ik} w_{kj}\right)^2$$

based on the notation of Problem 16. Now derive the block descent algorithm with multiplicative updates

$$v_{ik}^{n+1} = v_{ik}^n \frac{\sum_j u_{ij} w_{kj}^n}{\sum_j b_{ij}^n w_{kj}^n}$$

and

$$w_{kj}^{n+1} = w_{kj}^n \frac{\sum_i u_{ij} v_{ik}^n}{\sum_i b_{ij}^n v_{ik}^n}.$$

18. In Problem 16 calculate the partial derivative

$$\frac{\partial}{\partial v_{il}} L(\mathbf{V}, \mathbf{W}) = \sum_j w_{lj} \left(\frac{u_{ij}}{\sum_k v_{ik} w_{kj}} - 1 \right).$$

Show that the conditions $\min\{v_{il}, -\frac{\partial}{\partial v_{il}} L(\mathbf{V}, \mathbf{W})\} = 0$ for all pairs (i, l) are both necessary and sufficient for \mathbf{V} to maximize $L(\mathbf{V}, \mathbf{W})$ when \mathbf{W} is fixed. The same conditions apply in minimizing the criterion $\|\mathbf{U} - \mathbf{V}\mathbf{W}\|_F^2$ of Problem 17 with different partial derivatives.

19. In the matrix factorizations described in Problems 16 and 17, it may be worthwhile shrinking the estimates of the entries of \mathbf{V} and \mathbf{W} toward 0 [211]. Let λ and μ be positive constants, and consider the penalized objective functions

$$\begin{aligned} l(\mathbf{V}, \mathbf{W}) &= L(\mathbf{V}, \mathbf{W}) - \lambda \sum_i \sum_k v_{ik} - \mu \sum_k \sum_j w_{kj} \\ r(\mathbf{V}, \mathbf{W}) &= \|\mathbf{U} - \mathbf{V}\mathbf{W}\|_F^2 + \lambda \sum_i \sum_k v_{ik}^2 + \mu \sum_k \sum_j w_{kj}^2 \end{aligned}$$

with lasso and ridge penalties, respectively. Derive the block ascent updates

$$v_{ik}^{n+1} = v_{ik}^n \frac{\sum_j u_{ij} \frac{w_{kj}^n}{b_{ij}^n}}{\sum_j w_{kj}^n + \lambda}, \quad w_{kj}^{n+1} = w_{kj}^n \frac{\sum_i u_{ij} \frac{v_{ik}^n}{b_{ij}^n}}{\sum_i v_{ik}^n + \mu}$$

for $l(\mathbf{V}, \mathbf{W})$ and the block descent updates

$$v_{ik}^{n+1} = v_{ik}^n \frac{\sum_j u_{ij} w_{kj}^n}{\sum_j b_{ij}^n w_{kj}^n + \lambda v_{ik}^n}, \quad w_{kj}^{n+1} = w_{kj}^n \frac{\sum_i u_{ij} v_{ik}^n}{\sum_i b_{ij}^n v_{ik}^n + \mu w_{kj}^n}$$

for $r(\mathbf{V}, \mathbf{W})$. These updates maintain positivity. Shrinkage is obvious, with stronger shrinkage for the lasso penalty with small parameters.

20. Let $\mathbf{y}_1, \dots, \mathbf{y}_m$ be a random sample from a multivariate normal distribution on \mathbb{R}^p . Example 6.5.7 demonstrates that the sample mean $\bar{\mathbf{y}}$ and sample variance matrix \mathbf{S} are the maximum likelihood estimates of the theoretical mean $\boldsymbol{\mu}$ and variance $\boldsymbol{\Omega}$. The implicit assumption here is that $m \geq p$ and \mathbf{S} is invertible. Unfortunately, \mathbf{S} is singular whenever $m < p$. Furthermore, the entries of \mathbf{S} typically have high variance when $m \geq p$. To avoid these problems, Levina et al. [177] pursue lasso penalized estimation of $\boldsymbol{\Omega}^{-1}$. If we assume that $\boldsymbol{\Omega}$ is invertible and let $\boldsymbol{\Omega} = \mathbf{L}\mathbf{L}^*$ be its Cholesky decomposition, then $\boldsymbol{\Omega}^{-1} = (\mathbf{L}^*)^{-1}\mathbf{L}^{-1} = \mathbf{R}\mathbf{R}^*$ for the upper triangular matrix $\mathbf{R} = (r_{ij}) = (\mathbf{L}^*)^{-1}$. With the understanding $\hat{\boldsymbol{\mu}} = \bar{\mathbf{y}}$, show that the loglikelihood of the sample is

$$m \ln \det \mathbf{R} - \frac{m}{2} \operatorname{tr}(\mathbf{R}^* \mathbf{S} \mathbf{R}) = m \sum_i \ln r_{ii} - \frac{m}{2} \sum_j \mathbf{r}_j^* \mathbf{S} \mathbf{r}_j,$$

where \mathbf{r}_j is column j of \mathbf{R} . In lasso penalized estimation of \mathbf{R} , we minimize the objective function

$$f(\mathbf{R}) = -m \sum_i \ln r_{ii} + \frac{m}{2} \sum_j \mathbf{r}_j^* \mathbf{S} \mathbf{r}_j + \lambda \sum_{j>i} |r_{ij}|.$$

The diagonal entries of \mathbf{R} are not penalized because we want \mathbf{R} to be invertible. Why is $f(\mathbf{R})$ a convex function? For $r_{ij} \neq 0$, show that

$$\begin{aligned} \frac{\partial}{\partial r_{ij}} f(\mathbf{R}) &= -1_{\{j=i\}} \frac{m}{r_{ii}} + m s_{ii} r_{ij} + m \sum_{k \neq i} s_{ik} r_{kj} \\ &\quad + 1_{\{j \neq i\}} \begin{cases} \lambda & r_{ij} > 0 \\ -\lambda & r_{ij} < 0. \end{cases} \end{aligned}$$

Demonstrate that this leads to the cyclic coordinate descent update

$$\hat{r}_{ii} = \frac{-\sum_{k \neq i} s_{ik} r_{ki} + \sqrt{(\sum_{k \neq i} s_{ik} r_{ki})^2 + 4s_{ii}}}{2s_{ii}}.$$

Finally for $j \neq i$, demonstrate that the cyclic coordinate descent update chooses

$$\hat{r}_{ij} = -\frac{m \sum_{k \neq i} s_{ik} r_{kj} + \lambda}{m s_{ii}}$$

when this quantity is positive, it chooses

$$\hat{r}_{ij} = -\frac{m \sum_{k \neq i} s_{ik} r_{kj} - \lambda}{m s_{ii}}$$

when this second quantity is negative, and it chooses 0 otherwise. In organizing cyclic coordinate descent, it is helpful to retain and periodically update the sums $\sum_{k \neq i} s_{ik} r_{kj}$. The matrix \mathbf{R} can be traversed column by column.

14

Convex Calculus

14.1 Introduction

Two generations of mathematicians have labored to extend the machinery of differential calculus to convex functions. For many purposes it is convenient to generalize the definition of a convex function $f(\mathbf{x})$ to include the possibility that $f(\mathbf{x}) = \infty$. This maneuver has the advantage of allowing one to enlarge the domain of a convex function $f(\mathbf{x})$ defined on a convex set $C \subset \mathbb{R}^n$ to all of \mathbb{R}^n by the simple device of setting $f(\mathbf{x}) = \infty$ for $\mathbf{x} \notin C$. Many of the results for finite-valued convex functions generalize successfully in this setting. For instance, convex functions can still be characterized by their epigraphs and their satisfaction of Jensen's inequality.

The notion of the subdifferential $\partial f(\mathbf{x})$ of a convex function $f(\mathbf{x})$ has been a particularly fertile idea. This set consists of all vectors \mathbf{g} satisfying the supporting hyperplane inequality $f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^*(\mathbf{y} - \mathbf{x})$ for all \mathbf{y} . These vectors \mathbf{g} are called subgradients. If $f(\mathbf{x})$ is differentiable at \mathbf{x} , then $\partial f(\mathbf{x})$ reduces to the single vector $\nabla f(\mathbf{x})$. The subdifferential enjoys such familiar properties as

$$\begin{aligned}\partial[\alpha f(\mathbf{x})] &= \alpha \partial f(\mathbf{x}), \quad \alpha > 0 \\ \partial[f(\mathbf{x}) + g(\mathbf{x})] &= \partial f(\mathbf{x}) + \partial g(\mathbf{x}) \\ \partial[f \circ g(\mathbf{x})] &= dg(\mathbf{x})^* \partial f(\mathbf{y})|_{\mathbf{y}=g(\mathbf{x})}\end{aligned}$$

under the right hypotheses. Fermat's principle generalizes in the sense that \mathbf{y} furnishes a minimum of $f(\mathbf{x})$ if and only if $\mathbf{0} \in \partial f(\mathbf{y})$. A version

of the mean value theorem is true, and, properly interpreted, the Lagrange multiplier rule for a minimum remains valid. Perhaps more remarkable are the Fenchel conjugate and the formula for the subdifferential of the maximum of a finite collection of functions. The price of these successes is a theory more complicated than that encountered in classical calculus.

This chapter takes up the expository challenge of explaining these new concepts in the simplest possible terms. Fortunately, sacrificing generality for clarity does not mean losing sight of interesting applications. Convex calculus is an incredibly rich amalgam of ideas from analysis, linear algebra, and geometry. It has been instrumental in the construction of new algorithms for the solution of convex programs and their duals. Many readers will want to follow our brief account by pursuing the deeper treatises [13, 17, 131, 221, 226].

14.2 Notation

Although we allow ∞ for the value of a convex function, a function that is everywhere infinite is too boring to be of much interest. We will rule out such improper convex functions and the value $-\infty$ for a convex function. The convex set $\{\mathbf{x} \in \mathbb{R}^n : f(\mathbf{x}) < \infty\}$ is called the essential domain of $f(\mathbf{x})$ and abbreviated $\text{dom}(f)$. We denote the closure of a set C by $\text{cl } C$ and the convex hull of C by $\text{conv } C$.

The notion of lower semicontinuity introduced in Sect. 2.6 turns out to be crucial in many convexity arguments. Lower semicontinuity is equivalent to the epigraph $\text{epi}(f)$ being closed, and for this reason mathematicians call a lower semicontinuous function closed. Again, it makes sense to extend a finite closed function $f(\mathbf{x})$ with closed domain to all of \mathbb{R}^n by setting $f(\mathbf{x}) = \infty$ outside $\text{dom}(f)$. The fact that a closed convex function $f(\mathbf{x})$ has a closed convex epigraph permits application of the geometric separation property of convex sets described in Proposition 6.2.3. As an example of a closed convex function, consider

$$f(x) = \begin{cases} x \ln x - x & x > 0 \\ 0 & x = 0 \\ \infty & x < 0. \end{cases}$$

If on the one hand we redefine $f(0) < 0$, then $f(x)$ fails to be convex. If on the other hand we redefine $f(0) > 0$, then $f(x)$ fails to be closed.

14.3 Fenchel Conjugates

The Fenchel conjugate was defined in Example 1.2.6 for real-valued functions of a real argument. This transform, which is in some ways the

optimization analogue of the Fourier transform, has profound consequences in many branches of mathematics. The Fenchel conjugate generalizes to

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - f(\mathbf{x})] \quad (14.1)$$

for functions $f(\mathbf{x})$ mapping \mathbb{R}^n into $(-\infty, \infty]$. The conjugate $f^*(\mathbf{y})$ is always closed and convex even when $f(\mathbf{x})$ is neither. Condition (j) in the next proposition rules out improper conjugate functions.

On first contact definition (14.1) is frankly a little mysterious. So too is the definition of the Fourier transform. Readers are advised to exercise patience and suspend their initial skepticism for several reasons. The Fenchel conjugate encodes the solutions to a family of convex optimization problems. It is one of the keys to understanding convex duality and serves as a device for calculating subdifferentials. Finally, the Fenchel biconjugate $f^{**}(\mathbf{x})$ provides a practical way of convexifying $f(\mathbf{x})$. Indeed, the biconjugate has the geometric interpretation of falling below $f(\mathbf{x})$ and above any supporting hyperplane minorizing $f(\mathbf{x})$. This claim follows from our subsequent proof of the Fenchel-Moreau theorem and a double application of item (f) in the next proposition.

Proposition 14.3.1 *The Fenchel conjugate enjoys the following properties:*

- (a) *If $g(\mathbf{x}) = f(\mathbf{x} - \mathbf{v})$, then $g^*(\mathbf{y}) = f^*(\mathbf{y}) + \mathbf{v}^* \mathbf{y}$.*
- (b) *If $g(\mathbf{x}) = f(\mathbf{x}) - \mathbf{v}^* \mathbf{x} - c$, then $g^*(\mathbf{y}) = f^*(\mathbf{y} + \mathbf{v}) + c$.*
- (c) *If $g(\mathbf{x}) = \alpha f(\mathbf{x})$ for $\alpha > 0$, then $g^*(\mathbf{y}) = \alpha f^*(\alpha^{-1} \mathbf{y})$.*
- (d) *If $g(\mathbf{x}) = f(\mathbf{M}\mathbf{x})$, then $g^*(\mathbf{y}) = f^*[(\mathbf{M}^{-1})^* \mathbf{y}]$ for \mathbf{M} invertible.*
- (e) *If $g(\mathbf{x}) = \sum_{i=1}^n f_i(x_i)$ for $\mathbf{x} \in \mathbb{R}^n$, then $g^*(\mathbf{y}) = \sum_{i=1}^n f_i^*(y_i)$.*
- (f) *If $g(\mathbf{x}) \leq h(\mathbf{x})$ for all \mathbf{x} , then $g^*(\mathbf{y}) \geq h^*(\mathbf{y})$ for all \mathbf{y} .*
- (g) *For all \mathbf{x} and \mathbf{y} , $f^*(\mathbf{y}) + f(\mathbf{x}) \geq \mathbf{y}^* \mathbf{x}$.*
- (h) *The conjugate $f^*(\mathbf{y})$ is convex.*
- (i) *The conjugate $f^*(\mathbf{y})$ is closed.*
- (j) *If $f(\mathbf{x})$ satisfies $f(\mathbf{x}) \geq \mathbf{z}^* \mathbf{x} + c$ for all \mathbf{x} , then $f^*(\mathbf{z}) \leq -c$.*
- (k) *$f^*(\mathbf{0}) = -\inf_{\mathbf{x}} f(\mathbf{x})$.*

Proof: All of these claims are direct consequences of definition (14.1). For instance, claim (d) follows from

$$g^*(\mathbf{y}) = \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - f(\mathbf{M}\mathbf{x})] = \sup_{\mathbf{z}} [\mathbf{y}^* \mathbf{M}^{-1} \mathbf{z} - f(\mathbf{z})].$$

Claims (h) and (i) stem from the convexity and continuity of the affine functions $\mathbf{y} \mapsto \mathbf{y}^* \mathbf{x} - f(\mathbf{x})$ and the closure properties of convex and lower semicontinuous functions under suprema. Finally, claim (j) follows from the inequality $f^*(\mathbf{z}) \leq \sup_{\mathbf{x}} [(\mathbf{z} - \mathbf{z})^* \mathbf{x} - c] = -c$. ■

Example 14.3.1 *Conjugate of a Strictly Convex Quadratic*

For a well-behaved function $f(\mathbf{x})$, one can find $f^*(\mathbf{y})$ by setting the gradient $\mathbf{y} - \nabla f(\mathbf{x})$ equal to $\mathbf{0}$ and solving for \mathbf{x} . For example, consider the quadratic $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x}$ defined by the positive definite matrix \mathbf{A} . The gradient condition becomes $\mathbf{y} - \mathbf{A} \mathbf{x} = \mathbf{0}$ with solution $\mathbf{x} = \mathbf{A}^{-1} \mathbf{y}$. This result gives the conjugate $f^*(\mathbf{y}) = \frac{1}{2} \mathbf{y}^* \mathbf{A}^{-1} \mathbf{y}$. For the general convex quadratic

$$f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c,$$

rule (b) of Proposition 14.3.1 implies

$$f^*(\mathbf{y}) = \frac{1}{2} (\mathbf{y} - \mathbf{b})^* \mathbf{A}^{-1} (\mathbf{y} - \mathbf{b}) - c.$$

For instance, the univariate function $f(x) = \frac{1}{2} x^2$ is self-conjugate. According to rule (e) of Proposition 14.3.1, the multivariate function $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2$ is also self-conjugate. Rule (g) is called the Fenchel-Young inequality. In the current example it amounts to

$$\frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \frac{1}{2} \mathbf{y}^* \mathbf{A}^{-1} \mathbf{y} \geq \mathbf{y}^* \mathbf{x},$$

a surprising result in its own right. ■

Example 14.3.2 *Entropy as a Fenchel Conjugate*

Consider the convex function $f(\mathbf{x}) = \ln(\sum_{j=1}^n e^{x_j})$. If $f^*(\mathbf{y})$ is finite, then setting the gradient of $\mathbf{y}^* \mathbf{x} - f(\mathbf{x})$ with respect to \mathbf{x} equal to $\mathbf{0}$ entails

$$y_i = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}.$$

It follows that the y_i are positive and sum to 1. Furthermore,

$$f^*(\mathbf{y}) = \sum_{i=1}^n y_i \left[\ln y_i + \ln \left(\sum_{j=1}^n e^{x_j} \right) \right] - \ln \left(\sum_{j=1}^n e^{x_j} \right) = \sum_{i=1}^n y_i \ln y_i.$$

With the understanding that $0 \ln 0 = 0$, the entropy formula

$$f^*(\mathbf{y}) = \sum_{i=1}^n y_i \ln y_i$$

remains valid when \mathbf{y} occurs on the boundary of the unit simplex. Indeed, if \mathbf{y} belongs to the unit simplex but $y_i = 0$, then we send x_i to $-\infty$ and reduce calculation of the conjugate to the setting $\mathbf{y} \in \mathbb{R}^{n-1}$. If $y_i < 0$, then the sequence $\mathbf{x}_k = -k\mathbf{e}_i$ gives

$$\lim_{k \rightarrow \infty} \mathbf{y}^* \mathbf{x}_k - f(\mathbf{x}_k) = \lim_{k \rightarrow \infty} \ln \left(\frac{e^{-ky_i}}{n-1+e^{-k}} \right) = \infty.$$

Similarly, if all $y_i \geq 0$ but $\sum_{j=1}^n y_j \neq 1$, then one of the two sequences $\mathbf{x}_k = \pm k\mathbf{1}$ compels the same conclusion $f^*(\mathbf{y}) = \infty$. ■

Examples 1.2.6 and 14.3.1 obey the curious rule $f^{**}(\mathbf{x}) = f(\mathbf{x})$. This duality relation is more widely true. An affine function $f(\mathbf{x}) = \mathbf{z}^* \mathbf{x} + c$ provides another example. This fact follows readily from the form

$$f^*(\mathbf{y}) = \begin{cases} -c & \mathbf{y} = \mathbf{z} \\ \infty & \mathbf{y} \neq \mathbf{z} \end{cases}$$

of the conjugate. One cannot expect duality to hold for all functions because Proposition 14.3.1 requires $f^{**}(\mathbf{x})$ to be convex and closed. Remarkably, the combination of these two conditions is both necessary and sufficient for proper functions.

Proposition 14.3.2 (Fenchel-Moreau) *A proper function $f(\mathbf{x})$ from \mathbb{R}^n to $(-\infty, \infty]$ satisfies the duality relation $f^{**}(\mathbf{x}) = f(\mathbf{x})$ for all \mathbf{x} if and only if it is closed and convex.*

Proof: Suppose the duality relation holds. Being the conjugate of a conjugate, $f(\mathbf{x}) = f^{**}(\mathbf{x})$ is closed and convex. This proves that the stated conditions are necessary for duality.

To demonstrate the converse, first note that the Fenchel-Young inequality $f(\mathbf{x}) \geq \mathbf{y}^* \mathbf{x} - f^*(\mathbf{y})$ implies

$$f(\mathbf{x}) \geq \sup_{\mathbf{y}} [\mathbf{y}^* \mathbf{x} - f^*(\mathbf{y})] = f^{**}(\mathbf{x}). \tag{14.2}$$

In other words, the epigraph $\text{epi}(f^{**})$ of $f^{**}(\mathbf{x})$ contains the epigraph $\text{epi}(f)$ of $f(\mathbf{x})$. Verifying the reverse containment $\text{epi}(f^{**}) \subset \text{epi}(f)$ proves the converse of the proposition. Our general strategy for establishing containment is to exploit the separation properties of closed convex sets. As already mentioned, the convexity of $f(\mathbf{x})$ entails the convexity of $\text{epi}(f)$, and the lower semicontinuity of $f(\mathbf{x})$ entails the closedness of $\text{epi}(f)$.

We first show that $f(\mathbf{x})$ dominates some affine function and hence that the conjugate function $f^*(\mathbf{y})$ is proper. Suppose \mathbf{x}_0 satisfies $f(\mathbf{x}_0) < \infty$. Given that $\text{epi}(f)$ is closed and convex, we can separate it from the exterior point $[\mathbf{x}_0, f(\mathbf{x}_0) - 1]$ by a hyperplane. Thus, there exists a vector \mathbf{v} and scalars η and ν such that

$$\mathbf{v}^* \mathbf{x} + \eta r \geq \nu > \mathbf{v}^* \mathbf{x}_0 + \eta[f(\mathbf{x}_0) - 1] \tag{14.3}$$

for all \mathbf{x} and $r \geq f(\mathbf{x})$. Sending r to ∞ demonstrates that $\eta \geq 0$. Setting $\mathbf{x} = \mathbf{x}_0$ rules out the possibility $\eta = 0$. Finally, dividing inequality (14.3) by $\eta > 0$, replacing r by $f(\mathbf{x})$, and rearranging the result yield an affine function $\mathbf{z}^*\mathbf{x} + c$ positioned below $f(\mathbf{x})$.

Now suppose (\mathbf{y}, β) is in $\text{epi}(f^{**})$ but not in $\text{epi}(f)$. Proposition 6.2.3 guarantees the existence of a vector-scalar pair (\mathbf{v}, γ) and a constant $\epsilon > 0$ such that

$$\mathbf{v}^*\mathbf{y} + \gamma\beta \leq \mathbf{v}^*\mathbf{x} + \gamma\alpha - \epsilon$$

for all $(\mathbf{x}, \alpha) \in \text{epi}(f)$. Sending α to ∞ shows that $\gamma \geq 0$. If $\gamma > 0$, then

$$g(\mathbf{x}) = \gamma^{-1}\mathbf{v}^*(\mathbf{y} - \mathbf{x}) + \beta + \gamma^{-1}\epsilon \leq \alpha$$

for all $\alpha \geq f(\mathbf{x})$. Hence, $g(\mathbf{x})$ is an affine function positioned below $f(\mathbf{x})$, and a double application of part (f) of Proposition 14.3.1 implies

$$f^{**}(\mathbf{y}) \geq g^{**}(\mathbf{y}) = g(\mathbf{y}) > \beta,$$

contradicting the choice of $(\mathbf{y}, \beta) \in \text{epi}(f^{**})$.

Completing the proof now requires eliminating the possibility $\gamma = 0$. If we multiply the inequality $\mathbf{v}^*\mathbf{y} - \mathbf{v}^*\mathbf{x} + \epsilon \leq 0$ by $\delta > 0$ and add it to the previous inequality $\mathbf{z}^*\mathbf{x} + c \leq f(\mathbf{x})$, then we arrive at

$$h(\mathbf{x}) = \mathbf{z}^*\mathbf{x} + c + \delta\mathbf{v}^*(\mathbf{y} - \mathbf{x}) + \delta\epsilon \leq f(\mathbf{x}).$$

The conclusion

$$f^{**}(\mathbf{y}) \geq h^{**}(\mathbf{y}) = h(\mathbf{y}) = \mathbf{z}^*\mathbf{y} + c + \delta\epsilon$$

for all $\delta > 0$ can only be true if $f^{**}(\mathbf{y}) = \infty$, which is also incompatible with $(\mathbf{y}, \beta) \in \text{epi}(f^{**})$. ■

The left panel of Fig. 14.1 illustrates the relationship between a function $f(x)$ on the real line and its Fenchel conjugate $f^*(y)$. According to the Fenchel-Young inequality, the line with slope y and intercept $-f^*(y)$ falls below $f(x)$. The curve and the line intersect when $y = f'(x)$. The right panel of Fig. 14.1 shows that the biconjugate $f^{**}(y)$ is the greatest convex function lying below $f(x)$. The biconjugate is formed by taking the pointwise supremum of the supporting lines.

Example 14.3.3 *Perspective of a Function*

Let $f(\mathbf{x})$ be a closed convex function. The perspective of $f(\mathbf{x})$ is the function $g(\mathbf{x}, t) = tf(t^{-1}\mathbf{x})$ defined for $t > 0$. On this domain $g(\mathbf{x}, t)$ is closed and convex owing to the representation

$$g(\mathbf{x}, t) = t \sup_{\mathbf{y}} [t^{-1}\mathbf{x}^*\mathbf{y} - f^*(\mathbf{y})] = \sup_{\mathbf{y}} [\mathbf{x}^*\mathbf{y} - tf^*(\mathbf{y})]$$

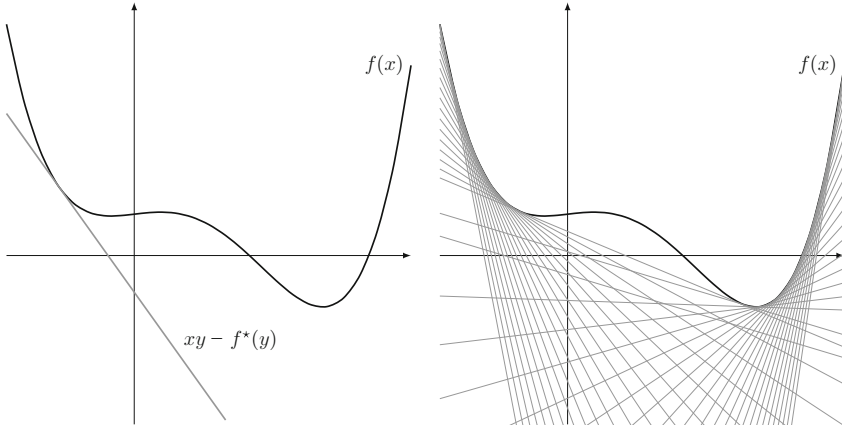


FIGURE 14.1. *Left panel:* The Fenchel-Young inequality for $f(x)$ and its conjugate $f^*(y)$. *Right panel:* The envelope of supporting lines defined by $f^*(y)$

and the linearity of the map $(\mathbf{x}, t) \mapsto \mathbf{x}^* \mathbf{y} - t f^*(\mathbf{y})$. For instance, the choice $f(\mathbf{x}) = \mathbf{x}^* \mathbf{M} \mathbf{x}$ for \mathbf{M} positive semidefinite yields the convexity of $t^{-1} \mathbf{x}^* \mathbf{M} \mathbf{x}$. The choice $f(x) = -\ln x$ shows that the relative entropy $g(x, t) = t \ln t - t \ln x$ is convex. Finally, the function

$$g(\mathbf{x}) = (\mathbf{c}^* \mathbf{x} + d) f\left(\frac{\mathbf{A} \mathbf{x} + \mathbf{b}}{\mathbf{c}^* \mathbf{x} + d}\right)$$

is closed and convex on the domain $\mathbf{c}^* \mathbf{x} + d > 0$ whenever the function $f(\mathbf{x})$ is closed and convex. ■

Example 14.3.4 *Indicator and Support Functions*

Every set C can be represented by its indicator function

$$\delta_C(\mathbf{x}) = \begin{cases} 0 & \mathbf{x} \in C \\ \infty & \mathbf{x} \notin C. \end{cases}$$

If C is closed and convex, then it is easy to check that $\delta_C(\mathbf{x})$ is a closed convex function. One reason for making the substitution of ∞ for 0 and 0 for 1 in defining an indicator is that it simplifies the Fenchel conjugate

$$\delta_C^*(\mathbf{y}) = \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - \delta_C(\mathbf{x})] = \sup_{\mathbf{x} \in C} \mathbf{y}^* \mathbf{x}.$$

The function $\delta_C^*(\mathbf{y})$ is called the support function of C . Proposition 14.3.2 implies that the Fenchel biconjugate $\delta_C^{**}(\mathbf{z})$ equals $\delta_C(\mathbf{z})$.

It turns out that support functions with full essential domains are the same as sublinear functions with full essential domains. A function $h(\mathbf{u})$ is said to be sublinear whenever

$$h(\alpha \mathbf{u} + \beta \mathbf{v}) \leq \alpha h(\mathbf{u}) + \beta h(\mathbf{v})$$

holds for all points \mathbf{u} and \mathbf{v} and nonnegative scalars α and β . Sublinearity is an amalgam of homogeneity, $h(\lambda\mathbf{v}) = \lambda h(\mathbf{v})$ for $\lambda > 0$, and convexity, $h[\lambda\mathbf{u} + (1 - \lambda)\mathbf{v}] \leq \lambda h(\mathbf{u}) + (1 - \lambda)h(\mathbf{v})$ for $\lambda \in [0, 1]$. One can easily check that a support function is sublinear. To prove the converse, we first demonstrate that the conjugate $h^*(\mathbf{y})$ of a sublinear function is an indicator function. Indeed, the identity

$$\begin{aligned}\lambda h^*(\mathbf{y}) &= \lambda \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - h(\mathbf{x})] \\ &= \sup_{\mathbf{x}} [\mathbf{y}^*(\lambda\mathbf{x}) - h(\lambda\mathbf{x})] \\ &= \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - h(\mathbf{x})] \\ &= h^*(\mathbf{y})\end{aligned}$$

compels $h^*(\mathbf{y})$ to equal 0 or ∞ . When $h(\mathbf{x})$ is finite-valued, Proposition 6.4.1 requires it to be continuous and consequently closed. Thus, Fenchel duality implies that $h(\mathbf{x})$ equals the support function of the closed convex set $C = \{\mathbf{y} \in \mathbb{R}^n : h^*(\mathbf{y}) = 0\}$.

The support function $h(\mathbf{u})$ of a closed convex set C is finite valued if and only if C is bounded. Boundedness is clearly sufficient to guarantee that $h(\mathbf{u})$ is finite valued. To show that boundedness is necessary as well, suppose C is unbounded. Then part (c) of Problem 5 of Chap. 6 says that C contains a ray $\{\mathbf{u} + t\mathbf{v} : t \in [0, \infty)\}$. This forces $h(\mathbf{v}) = \sup_{\mathbf{x} \in C} \mathbf{v}^* \mathbf{x}$ to be infinite. ■

Example 14.3.5 Vector Dual Norms

If C equals the closed unit ball $B = \{\mathbf{x} : \|\mathbf{x}\|_{\dagger} \leq 1\}$ associated with a norm $\|\mathbf{x}\|_{\dagger}$ on \mathbb{R}^n , then the support function $\|\mathbf{y}\|_{\star} = \delta_B^*(\mathbf{y}) = \sup_{\mathbf{x} \in B} \mathbf{y}^* \mathbf{x}$ also qualifies as a norm. Verification of the norm properties for the dual norm $\|\mathbf{y}\|_{\star}$ is left to the reader as Problem 13. For a sublinear function such as $\|\mathbf{y}\|_{\star}$, the only things to check are that whether it is nonnegative and vanishes if and only if $\mathbf{x} = \mathbf{0}$. In defining $\|\mathbf{y}\|_{\star}$ one can clearly confine \mathbf{x} to the boundary of B . The generalized Cauchy-Schwarz inequality

$$\mathbf{y}^* \mathbf{x} \leq \|\mathbf{y}\|_{\star} \|\mathbf{x}\|_{\dagger} \tag{14.4}$$

follows directly from this observation. Equality in the generalized Cauchy-Schwarz inequality is attained for some \mathbf{x} on the boundary of B because the maximum of a continuous function over a compact set is attained.

There are many concrete examples of norms and their duals. The dual norm of $\|\mathbf{x}\|_1$ is $\|\mathbf{x}\|_{\infty}$ and vice versa. For $p^{-1} + q^{-1} = 1$, Example 6.6.3 and Problem 21 of Chap. 5 show that the norms $\|\mathbf{x}\|_p$ and $\|\mathbf{y}\|_q$ constitute another dual pair, with the generalized Cauchy-Schwarz inequality reducing to Hölder's inequality. Of course, the Euclidean norm $\|\mathbf{x}\| = \|\mathbf{x}\|_2$ is its own dual.

These examples are not accidents. The primary reason for calling $\|\mathbf{y}\|_*$ a dual norm is that taking the dual of the dual gives us back the original norm $\|\mathbf{x}\|_{\dagger}$. One way of deducing this fact is to construct the Fenchel conjugate of the convex function $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_{\dagger}^2$. In view of the generalized Cauchy-Schwarz inequality and the inequality $(\|\mathbf{x}\|_{\dagger} - \|\mathbf{y}\|_*)^2 \geq 0$, we have

$$\mathbf{y}^* \mathbf{x} - \frac{1}{2}\|\mathbf{x}\|_{\dagger}^2 \leq \|\mathbf{y}\|_* \|\mathbf{x}\|_{\dagger} - \frac{1}{2}\|\mathbf{x}\|_{\dagger}^2 \leq \frac{1}{2}\|\mathbf{y}\|_*^2.$$

On the other hand, suppose we choose a vector \mathbf{z} with $\|\mathbf{z}\|_{\dagger} = 1$ such that equality is attained in the generalized Cauchy-Schwarz inequality. Then for any scalar $s > 0$, we find

$$\mathbf{y}^*(s\mathbf{z}) - \frac{1}{2}\|s\mathbf{z}\|_{\dagger}^2 = s\|\mathbf{y}\|_* \|\mathbf{z}\|_{\dagger} - \frac{s^2}{2}\|\mathbf{z}\|_{\dagger}^2.$$

If we take $s = \|\mathbf{y}\|_*$, then this equality gives

$$\mathbf{y}^*(s\mathbf{z}) - \frac{1}{2}\|s\mathbf{z}\|_{\dagger}^2 = \frac{1}{2}\|\mathbf{y}\|_*^2.$$

In other words, $f^*(\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|_*^2$. Taking the conjugate of $f^*(\mathbf{y}) = \frac{1}{2}\|\mathbf{y}\|_*^2$ yields $f^{**}(\mathbf{x}) = f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|_{\dagger}^2$. Thus, the original norm $\|\mathbf{x}\|_{\dagger}$ is dual to the dual norm $\|\mathbf{y}\|_*$. ■

Example 14.3.6 *Matrix Dual Norms*

One can also define dual norms of matrix norms using the Frobenius inner product $\langle \mathbf{Y}, \mathbf{X} \rangle = \text{tr}(\mathbf{Y}^* \mathbf{X})$. Under this matrix inner product, the Frobenius matrix norm $\|\mathbf{X}\|_F$ is self-dual. The easiest way to deduce this fact is to observe that the Frobenius norm can be calculated by stacking the columns of \mathbf{X} to form a vector and then taking the Euclidean norm of the vector. Column stacking is clearly compatible with the exchange of the Frobenius inner product for the Euclidean inner product.

Under the Frobenius inner product, calculation of the dual norm of the matrix spectral norm is more subtle. The most illuminating approach takes a detour through Fan’s inequality and the singular value decomposition (svd) covered in Appendices A.4 and A.5. Any matrix \mathbf{X} has an svd representation $\mathbf{P}\Sigma\mathbf{Q}^*$, where \mathbf{P} and \mathbf{Q} are orthogonal matrices and Σ is a diagonal matrix with nonnegative entries σ_i arranged in decreasing order along its diagonal. The columns of \mathbf{P} and \mathbf{Q} are referred to as singular vectors and the diagonal entries of Σ as singular values. The svd immediately yields the spectral decompositions $\mathbf{X}\mathbf{X}^* = \mathbf{P}\Sigma^2\mathbf{P}^*$ and $\mathbf{X}^*\mathbf{X} = \mathbf{Q}\Sigma^2\mathbf{Q}^*$ and consequently the spectral norm $\|\mathbf{X}\| = \sigma_1$. If \mathbf{Y} has svd $\mathbf{R}\Omega\mathbf{S}^*$ with $\Omega = \text{diag}(\omega_i)$, then equality is attained in Fan’s inequality

$$\text{tr}(\mathbf{Y}^* \mathbf{X}) \leq \sum_i \omega_i \sigma_i$$

when $\mathbf{R} = \mathbf{P}$ and $\mathbf{S} = \mathbf{Q}$. The matrix \mathbf{X} with $\|\mathbf{X}\| = \sigma_1 \leq 1$ giving the maximum value of $\text{tr}(\mathbf{Y}^* \mathbf{X})$ has the same singular vectors as \mathbf{Y} and the singular values $\sigma_i = 1$ for $\omega_i > 0$. It follows that the dual norm of the spectral norm equals

$$\|\mathbf{Y}\|_* = \sum_i \omega_i.$$

This dual norm is also called the nuclear norm or the trace norm. ■

Example 14.3.7 *Cones and Polar Cones*

If C is a cone, then the Fenchel conjugate of its indicator function $\delta_C(\mathbf{x})$ turns out to be the indicator function of the polar cone

$$C^\circ = \{\mathbf{y} : \mathbf{y}^* \mathbf{x} \leq 0, \forall \mathbf{x} \in C\}.$$

This assertion follows from the limits

$$\begin{aligned} \lim_{c \downarrow 0} \mathbf{y}^*(c\mathbf{x}) &= 0 \\ \lim_{c \uparrow \infty} \mathbf{y}^*(c\mathbf{x}) &= \begin{cases} \infty & \mathbf{y}^* \mathbf{x} > 0 \\ 0 & \mathbf{y}^* \mathbf{x} = 0 \\ -\infty & \mathbf{y}^* \mathbf{x} < 0 \end{cases} \end{aligned}$$

for any $\mathbf{x} \in C$. Although C may be neither convex nor closed, its polar C° is always both. The duality relation $\delta_{C^\circ}^{**}(\mathbf{x}) = \delta_C(\mathbf{x})$ for a closed convex cone C is equivalent to the set relation $C^{\circ\circ} = C$. Notice the analogy here to the duality relation $S^{\perp\perp} = S$ for subspaces under the orthogonal complement operator \perp .

As a concrete example, let us calculate the polar cone of the set S_+^n of $n \times n$ positive semidefinite matrices under the Frobenius inner product $\langle \mathbf{A}, \mathbf{B} \rangle = \text{tr}(\mathbf{A}\mathbf{B})$. If $\mathbf{A} \in S_+^n$ has eigenvalues $\lambda_1, \dots, \lambda_n$ with corresponding unit eigenvectors $\mathbf{u}_1, \dots, \mathbf{u}_n$, then

$$\text{tr}(\mathbf{A}\mathbf{B}) = \text{tr}\left(\sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^* \mathbf{B}\right) = \sum_{i=1}^n \lambda_i \mathbf{u}_i^* \mathbf{B} \mathbf{u}_i.$$

If $\mathbf{B} \in -S_+^n$, then it is clear that $\text{tr}(\mathbf{A}\mathbf{B}) \leq 0$. Thus, the polar cone contains $-S_+^n$. Conversely, suppose \mathbf{B} is in the polar cone of S_+^n , and choose $\mathbf{A} = \mathbf{v}\mathbf{v}^*$ for some nontrivial vector \mathbf{v} . The inequality $\text{tr}(\mathbf{v}\mathbf{v}^* \mathbf{B}) = \mathbf{v}^* \mathbf{B} \mathbf{v} \leq 0$ for all such \mathbf{v} implies that $\mathbf{B} \in -S_+^n$. Thus, the polar cone of S_+^n equals $-S_+^n$. ■

Example 14.3.8 *Log Determinant*

Minimization of the objective function (4.18) featured in Example 4.7.6 can be rephrased as calculating the Fenchel conjugate

$$g^*(\mathbf{N}) = \sup_{\mathbf{P}} [\text{tr}(\mathbf{N}\mathbf{P}) + c \ln \det \mathbf{P}]$$

for c positive and \mathbf{P} an $n \times n$ positive definite matrix. The essential domain of $g^*(\mathbf{N})$ is the set of negative definite matrices. Indeed, if \mathbf{N} falls outside this domain, then it possesses a unit eigenvector \mathbf{v} with nonnegative eigenvalue λ . The choice $\mathbf{P} = \mathbf{I} + s\mathbf{v}\mathbf{v}^*$ has eigenvalues 1 (multiplicity $n - 1$) and $1 + s$ (multiplicity 1). For $s > 0$ we calculate

$$\begin{aligned} \operatorname{tr}(\mathbf{N}\mathbf{P}) + c \ln \det \mathbf{P} &= \operatorname{tr} \mathbf{N} + s\lambda + c \ln \det(\mathbf{I} + s\mathbf{v}\mathbf{v}^*) \\ &= \operatorname{tr} \mathbf{N} + s\lambda + c \ln(1 + s), \end{aligned}$$

which tends to ∞ as s tends to ∞ . When \mathbf{N} is negative definite, our previous calculations gave the gradient $\mathbf{N} + c\mathbf{P}^{-1}$ of $\operatorname{tr}(\mathbf{N}\mathbf{P}) + c \ln \det \mathbf{P}$. Setting the gradient to $\mathbf{0}$ yields $\mathbf{P} = -c\mathbf{N}^{-1}$ and $g^*(\mathbf{N}) = -cn - c \ln \det[-c^{-1}\mathbf{N}]$. Because a Fenchel conjugate is convex, this line of argument establishes the log-concavity of $\det \mathbf{P}$ for \mathbf{P} positive definite. Example 6.3.12 presents a different proof of this fact. ■

14.4 Subdifferentials

Convex calculus revolves around the ideas of forward directional derivatives and supporting hyperplanes. At this juncture the reader may want to review Sect. 6.4 on the former topic. Appendix A.6 develops the idea of a semidifferential, the single most fruitful generalization of forward directional derivatives to date. For the sake of brevity henceforth, we will drop the adjective forward and refer to forward directional derivatives simply as directional derivatives.

Consider the absolute value function $f(x) = |x|$. At the point $x = 0$, the derivative $f'(x)$ does not exist. However, the directional derivatives $d_v f(0) = |v|$ are all well defined. Furthermore, the supporting hyperplane inequality $f(x) \geq f(0) + gx$ is valid for all x and all g with $|g| \leq 1$. The set $\partial f(0) = \{g : |g| \leq 1\}$ is called the subdifferential of $f(x)$ at $x = 0$. The notion of subdifferential is hardly limited to functions defined on the real line. Consider a convex function $f(\mathbf{x}) : \mathbb{R}^n \mapsto (-\infty, \infty]$. By convention the subdifferential $\partial f(\mathbf{x})$ is empty for $\mathbf{x} \notin \operatorname{dom}(f)$. For $\mathbf{x} \in \operatorname{dom}(f)$, the subdifferential $\partial f(\mathbf{x})$ is the set of vectors \mathbf{g} in \mathbb{R}^n such that the supporting hyperplane inequality

$$f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{g}^*(\mathbf{y} - \mathbf{x})$$

is valid for all \mathbf{y} . An element \mathbf{g} of the subdifferential is called a subgradient. If $f(\mathbf{x})$ is differentiable at \mathbf{x} , then as we prove later, its subdifferential collapses to its gradient.

For a nontrivial example, consider $f(\mathbf{x}) = \max\{x_1, x_2\}$ with domain \mathbb{R}^2 . Off the diagonal $D = \{\mathbf{x} \in \mathbb{R}^2 : x_1 = x_2\}$, the convex function $f(\mathbf{x})$ is differentiable. On D the subdifferential $\partial f(\mathbf{x})$ equals the unit simplex

$S = \{\mathbf{w} \in \mathbb{R}^2 : w_1 + w_2 = 1, w_1 \geq 0, w_2 \geq 0\}$. Indeed for points $\mathbf{x} \in D$ and $\mathbf{w} \in S$, one can easily demonstrate that

$$\max\{y_1, y_2\} \geq \max\{x_1, x_2\} + w_1(y_1 - x_1) + w_2(y_2 - x_2)$$

is valid for all \mathbf{y} by closely examining the two extreme cases $\mathbf{w} = (1, 0)^*$ and $\mathbf{w} = (0, 1)^*$. Conversely, we will prove later that any point \mathbf{w} outside S fails the supporting hyperplane test for some \mathbf{y} . The directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ equals $(1, 0)\mathbf{v} = v_1$ for $x_1 > x_2$ and $(0, 1)\mathbf{v} = v_2$ for $x_2 > x_1$. Example 4.4.4 shows that the directional derivative is $d_{\mathbf{v}}f(\mathbf{x}) = \max\{v_1, v_2\}$ on D . It is no accident that

$$d_{\mathbf{v}}f(\mathbf{x}) = \max\{\mathbf{w}^*\mathbf{v} : \mathbf{w} \in \partial f(\mathbf{x}) = S\} \quad (14.5)$$

on D . Indeed, the relationship (14.5) is generally true.

As a prelude to proving this fact, let us study the directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ more thoroughly. For $f(\mathbf{x})$ convex and \mathbf{x} an interior point of $\text{dom}(f)$, one can argue that $d_{\mathbf{v}}f(\mathbf{x})$ exists and is finite for all \mathbf{v} because any line segment starting at \mathbf{x} can be extended backwards and remain in $\text{dom}(f)$. Given that \mathbf{x} is internal to the segment, the results of Sect. 6.4 apply. In addition, one can show that $d_{\mathbf{v}}f(\mathbf{x})$ is sublinear in its argument \mathbf{v} . See Example 14.3.4 for the definition of sublinearity. A directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ is obviously homogeneous in \mathbf{v} . Because

$$d_{\mathbf{v}}f(\mathbf{x}) = \lim_{t \downarrow 0} \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}$$

is a pointwise limit of its difference quotients, and these are convex functions of \mathbf{v} , $d_{\mathbf{v}}f(\mathbf{x})$ is also convex. The monotonicity of the difference quotient implies that

$$f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x}) \geq td_{\mathbf{v}}f(\mathbf{x}).$$

Any vector \mathbf{g} satisfying $\mathbf{g}^*\mathbf{v} \leq d_{\mathbf{v}}f(\mathbf{x})$ for all \mathbf{v} therefore acts as a subgradient. Conversely, any subgradient \mathbf{g} must satisfy $\mathbf{g}^*\mathbf{v} \leq d_{\mathbf{v}}f(\mathbf{x})$ for all \mathbf{v} . Hence, we have the relation

$$\partial f(\mathbf{x}) = \{\mathbf{g} : \mathbf{g}^*\mathbf{v} \leq d_{\mathbf{v}}f(\mathbf{x}) \text{ for all } \mathbf{v}\}. \quad (14.6)$$

Despite this interesting logical equivalence, the question remains of whether any subgradient \mathbf{g} satisfies $\mathbf{g}^*\mathbf{v} = d_{\mathbf{v}}f(\mathbf{x})$ for a particular direction \mathbf{v} .

In constructing a subgradient \mathbf{g} with $\mathbf{g}^*\mathbf{v} = d_{\mathbf{v}}f(\mathbf{x})$, we view the map $\mathbf{u} \mapsto \mathbf{g}^*\mathbf{u}$ as a linear functional $\ell(\mathbf{u})$. Initially $\ell(\mathbf{u})$ is defined for the vector \mathbf{v} of interest by the requirement $\ell(\mathbf{v}) = d_{\mathbf{v}}f(\mathbf{x})$. We now show how to extend $\ell(\mathbf{u})$ to all of \mathbb{R}^n . By homogeneity,

$$\ell(\lambda\mathbf{v}) = \lambda\ell(\mathbf{v}) = d_{\lambda\mathbf{v}}f(\mathbf{x})$$

for $\lambda \geq 0$. This defines $\ell(\mathbf{u})$ on the ray $\{\lambda \mathbf{v} : \lambda \geq 0\}$. For $\lambda < 0$ we continue to define $\ell(\lambda \mathbf{v})$ by $\lambda \ell(\mathbf{v})$. Jensen's inequality

$$0 = d_{\mathbf{0}}f(\mathbf{x}) \leq \frac{1}{2}d_{-\mathbf{v}}f(\mathbf{x}) + \frac{1}{2}d_{\mathbf{v}}f(\mathbf{x})$$

then implies that

$$\ell(-\mathbf{v}) = -d_{\mathbf{v}}f(\mathbf{x}) \leq d_{-\mathbf{v}}f(\mathbf{x})$$

and assures us that $\ell(\mathbf{u})$ is dominated by $d_{\mathbf{u}}f(\mathbf{x})$ on the 1-dimensional subspace $\{\lambda \mathbf{v} : \lambda \in \mathbb{R}\}$. The remainder of the proof is supplied by the finite-dimensional version of the Hahn-Banach theorem.

Proposition 14.4.1 (Hahn-Banach) *Suppose the linear function $\ell(\mathbf{v})$ is defined on a subspace S of \mathbb{R}^n and dominated there by the sublinear function $h(\mathbf{v})$ defined on all of \mathbb{R}^n . Then $\ell(\mathbf{v})$ can be extended to a linear function that is dominated throughout \mathbb{R}^n by $h(\mathbf{v})$.*

Proof: The proof proceeds by induction on the dimension of S . Let \mathbf{u} be any point in \mathbb{R}^n outside S . It suffices to show that $\ell(\mathbf{v})$ can be consistently defined on the subspace T spanned by S and \mathbf{u} . For $\mathbf{v} \in S$ linearity requires

$$\ell(\mathbf{v} + \lambda \mathbf{u}) = \ell(\mathbf{v}) + \lambda \ell(\mathbf{u}),$$

and the crux of the matter is properly defining $\ell(\mathbf{u})$. For two points \mathbf{v} and \mathbf{w} in S , we have

$$\ell(\mathbf{v}) + \ell(\mathbf{w}) = \ell(\mathbf{v} + \mathbf{w}) \leq h(\mathbf{v} + \mathbf{w}) \leq h(\mathbf{v} - \mathbf{u}) + h(\mathbf{w} + \mathbf{u}).$$

It follows that

$$\ell(\mathbf{v}) - h(\mathbf{v} - \mathbf{u}) \leq h(\mathbf{w} + \mathbf{u}) - \ell(\mathbf{w}).$$

Because \mathbf{v} and \mathbf{w} are arbitrary, the left-hand side of this inequality is bounded above for \mathbf{u} fixed and the right-hand side is bounded below for \mathbf{u} fixed. The idea now is to define $\ell(\mathbf{u})$ to be any number α satisfying

$$\sup_{\mathbf{v} \in S} [\ell(\mathbf{v}) - h(\mathbf{v} - \mathbf{u})] \leq \alpha \leq \inf_{\mathbf{w} \in S} [h(\mathbf{w} + \mathbf{u}) - \ell(\mathbf{w})].$$

If $\lambda > 0$, then our choice of α entails

$$\begin{aligned} \ell(\mathbf{v} + \lambda \mathbf{u}) &= \lambda [\alpha + \ell(\lambda^{-1} \mathbf{v})] \\ &\leq \lambda [h(\lambda^{-1} \mathbf{v} + \mathbf{u}) - \ell(\lambda^{-1} \mathbf{v}) + \ell(\lambda^{-1} \mathbf{v})] \\ &= h(\mathbf{v} + \lambda \mathbf{u}). \end{aligned}$$

Similarly if $\lambda < 0$, then our choice of α entails

$$\begin{aligned} \ell(\mathbf{v} + \lambda \mathbf{u}) &= -\lambda [-\alpha + \ell(-\lambda^{-1} \mathbf{v})] \\ &\leq -\lambda [h(-\lambda^{-1} \mathbf{v} - \mathbf{u}) - \ell(-\lambda^{-1} \mathbf{v}) + \ell(-\lambda^{-1} \mathbf{v})] \\ &= h(\mathbf{v} + \lambda \mathbf{u}). \end{aligned}$$

Thus, $\ell(\mathbf{x})$ is dominated by $h(\mathbf{x})$ on T . ■

The next proposition summarizes the previous discussion and collects some further pertinent facts about subdifferentials. Recall that the notion of a support function mentioned in the proposition was defined in Example 14.3.4.

Proposition 14.4.2 *The subdifferential $\partial f(\mathbf{x})$ of a convex function $f(\mathbf{x})$ is a closed convex set for all \mathbf{x} . If \mathbf{x} is an interior point of $\text{dom}(f)$, then $\partial f(\mathbf{x})$ is nonempty and compact, the directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ exists and is finite for all \mathbf{v} , and $d_{\mathbf{v}}f(\mathbf{x})$ is the support function of $\partial f(\mathbf{x})$. If $f(\mathbf{x})$ is differentiable at \mathbf{x} , then $\partial f(\mathbf{x})$ coincides with the singleton set $\{\nabla f(\mathbf{x})\}$.*

Proof: The supporting hyperplane inequality is preserved under limits and convex combinations of subgradients. Hence, $\partial f(\mathbf{x})$ is closed and convex. For an interior point \mathbf{x} , the Hahn-Banach theorem shows that $\partial f(\mathbf{x})$ is nonempty and that

$$d_{\mathbf{v}}f(\mathbf{x}) = \max_{\mathbf{g} \in \partial f(\mathbf{x})} \mathbf{g}^* \mathbf{v},$$

so by definition $d_{\mathbf{v}}f(\mathbf{x})$ is the support function of $\partial f(\mathbf{x})$. To demonstrate that $\partial f(\mathbf{x})$ is compact, it suffices to show that it is bounded. To reach a contradiction, assume that $\mathbf{g}_m \in \partial f(\mathbf{x})$ satisfies $\lim_{m \rightarrow \infty} \|\mathbf{g}_m\| = \infty$. By passing to a subsequence if necessary, one can further assume that the sequence of unit vectors $\mathbf{v}_m = \|\mathbf{g}_m\|^{-1} \mathbf{g}_m$ converges to a unit vector \mathbf{v} . If the ball of radius ϵ centered at \mathbf{x} lies wholly within $\text{dom}(f)$, then the inequality

$$f(\mathbf{x} + \epsilon \mathbf{v}_m) \geq f(\mathbf{x}) + \epsilon \mathbf{g}_m^* \mathbf{v}_m = f(\mathbf{x}) + \epsilon \|\mathbf{g}_m\|$$

contradicts the boundedness of $f(\mathbf{y})$ within the ball guaranteed by Proposition 6.4.1. Finally, suppose $f(\mathbf{x})$ is differentiable at \mathbf{x} . For $\mathbf{g} \in \partial f(\mathbf{x})$ we have

$$f(\mathbf{x}) + \alpha \mathbf{g}^* \mathbf{v} \leq f(\mathbf{x} + \alpha \mathbf{v}) = f(\mathbf{x}) + \alpha df(\mathbf{x})\mathbf{v} + o(|\alpha|).$$

If we let $\mathbf{v} = \nabla f(\mathbf{x}) - \mathbf{g}$, then this inequality implies

$$0 \leq \alpha [\nabla f(\mathbf{x}) - \mathbf{g}]^* [\nabla f(\mathbf{x}) - \mathbf{g}] + o(|\alpha|) = \alpha \|\nabla f(\mathbf{x}) - \mathbf{g}\|^2 + o(|\alpha|).$$

Taking $\alpha < 0$ small now yields a contradiction unless $\nabla f(\mathbf{x}) - \mathbf{g} = \mathbf{0}$. ■

For a convex function $f(x)$ with domain the real line, the subdifferential $\partial f(x)$ equals the interval $[-d_{-1}f(x), d_1f(x)]$. For instance, the function $f(x) = |y - x|$ has

$$\partial f(x) = \begin{cases} -1 & x < y \\ [-1, 1] & x = y \\ 1 & x > y. \end{cases}$$

The behavior of $\partial f(x)$ at boundary points of $\text{dom}(f)$ is more erratic than the behavior at interior points. For example, the choice

$$f(x) = \begin{cases} -\sqrt{x} & x \in [0, 1] \\ \infty & \text{otherwise} \end{cases}$$

has subdifferential

$$\partial f(x) = \begin{cases} -\frac{1}{2\sqrt{x}} & x \in (0, 1) \\ [-1/2, \infty) & x = 1 \\ \emptyset & x \leq 0 \text{ or } x > 1. \end{cases}$$

Thus, at one boundary point of $\text{dom}(f)$ the subdifferential $\partial f(x)$ is empty, and at the other it is unbounded.

Here is the convex generalization of Fermat's stationarity condition.

Proposition 14.4.3 *A convex function $f(\mathbf{y})$ possesses a minimum at the point \mathbf{x} if and only if $\mathbf{0} \in \partial f(\mathbf{x})$.*

Proof: The inequality $f(\mathbf{y}) \geq f(\mathbf{x})$ for all \mathbf{y} is trivially equivalent to the inequality $f(\mathbf{y}) \geq f(\mathbf{x}) + \mathbf{0}^*(\mathbf{y} - \mathbf{x})$ for all \mathbf{y} . ■

The next proposition highlights the importance of the Fenchel conjugate in calculating subdifferentials.

Proposition 14.4.4 *For a convex function $f(\mathbf{x})$ with conjugate $f^*(\mathbf{y})$, assertions (a) and (b) from the following list*

(a) $f(\mathbf{x}) + f^*(\mathbf{y}) = \mathbf{y}^*\mathbf{x}$

(b) $\mathbf{y} \in \partial f(\mathbf{x})$

(c) $\mathbf{x} \in \partial f^*(\mathbf{y})$

are logically equivalent for a vector pair (\mathbf{x}, \mathbf{y}) . If $f(\mathbf{x})$ is closed, then both assertions are logically equivalent to assertion (c). Furthermore, the set of minima of $f(\mathbf{x})$ coincides with $\partial f^(\mathbf{0})$.*

Proof: A pair (\mathbf{x}, \mathbf{y}) satisfies condition (a) if and only if $\mathbf{y}^*\mathbf{x} - f(\mathbf{x})$ attains its maximum at \mathbf{x} for \mathbf{y} fixed. The latter condition is equivalent to the convex function $h(\mathbf{x}) = f(\mathbf{x}) - \mathbf{y}^*\mathbf{x}$ attaining its minimum. A brief calculation shows that $\partial h(\mathbf{x}) = \partial f(\mathbf{x}) - \mathbf{y}$. Proposition 14.4.3 therefore implies that condition (a) is equivalent to $\mathbf{0} \in \partial f(\mathbf{x}) - \mathbf{y}$, which is just a restatement of condition (b). When $f(\mathbf{x})$ is closed, $f^{**}(\mathbf{x}) = f(\mathbf{x})$, and we can reverse the roles of $f(\mathbf{x})$ and $f^*(\mathbf{y})$ and deduce the equivalence of condition (c). The final assertion of the proposition follows from the observation that $\mathbf{0} \in \partial f(\mathbf{x})$ if and only if $\mathbf{x} \in \partial f^*(\mathbf{0})$. ■

Example 14.4.1 *Subdifferential of the Indicator of a Closed Convex Set*

Let $\delta_C(\mathbf{x})$ be the indicator of a closed convex set C . Outside C the subdifferential $\partial\delta_C(\mathbf{x})$ is empty. For $\mathbf{x} \in C$ the calculation

$$\begin{aligned}\partial\delta_C(\mathbf{x}) &= \{\mathbf{y} : \delta_C(\mathbf{x}) + \delta_C^*(\mathbf{y}) = \mathbf{y}^*\mathbf{x}\} \\ &= \{\mathbf{y} : \sup_{\mathbf{z} \in C} \mathbf{y}^*\mathbf{z} = \mathbf{y}^*\mathbf{x}\} \\ &= \{\mathbf{y} : \mathbf{y}^*(\mathbf{z} - \mathbf{x}) \leq 0 \text{ for all } \mathbf{z} \in C\}\end{aligned}$$

identifies the subdifferential as the polar cone to the translated set $C - \mathbf{x}$. This polar cone is denoted $N_C(\mathbf{x})$ and called the normal cone to C at \mathbf{x} . The double polar cone $N_C(\mathbf{x})^\circ$ is termed the tangent cone to C at \mathbf{x} ; it is the smallest closed convex cone containing $C - \mathbf{x}$. Various special cases of $N_C(\mathbf{x})$ readily come to mind. For example, if \mathbf{x} belongs to the interior of C , then $N_C(\mathbf{x}) = \{\mathbf{0}\}$. Alternatively, if C is an affine subspace $S + \mathbf{x}$, then $N_C(\mathbf{x}) = S^\perp$, the orthogonal complement of the subspace S . When C is the affine subspace $\{\mathbf{x} : \mathbf{A}\mathbf{x} = \mathbf{b}\}$ corresponding to solutions of the linear equation $\mathbf{A}\mathbf{x} = \mathbf{b}$, then S equals the kernel of \mathbf{A} , and the fundamental theorem of linear algebra implies that S^\perp equals the range of the transpose \mathbf{A}^* [248]. ■

Example 14.4.2 *Subdifferential of a Norm*

Let B be the closed unit ball associated with a norm $\|\mathbf{x}\|_\dagger$ on \mathbb{R}^n . The dual norm $\|\mathbf{y}\|_\star$ defined in Example 14.3.5 coincides with the support function $\delta_B^*(\mathbf{y})$. These considerations imply

$$\begin{aligned}\partial\|\mathbf{y}\|_\star &= \{\mathbf{x} : \delta_B(\mathbf{x}) + \delta_B^*(\mathbf{y}) = \mathbf{y}^*\mathbf{x}\} \\ &= \{\mathbf{x} \in B : \|\mathbf{y}\|_\star = \mathbf{y}^*\mathbf{x}\}.\end{aligned}$$

The dual relation

$$\partial\|\mathbf{x}\|_\dagger = \{\mathbf{y} \in U : \|\mathbf{x}\|_\dagger = \mathbf{y}^*\mathbf{x}\}$$

holds for U the closed unit ball associated with $\|\mathbf{y}\|_\star$. In the case of the Euclidean norm, which is its own dual, the subdifferential coincides with the gradient $\|\mathbf{x}\|^{-1}\mathbf{x}$ when $\mathbf{x} \neq \mathbf{0}$. At the origin $\partial\|\mathbf{0}\| = B$. In general for any norm, $\mathbf{0} \in \partial\|\mathbf{x}\|_\dagger$ if and only if $\mathbf{x} = \mathbf{0}$. This is just a manifestation of the fact that $\mathbf{x} = \mathbf{0}$ is the unique minimum point of the norm. ■

Example 14.4.3 *Subdifferential of the Distance to a Closed Convex Set*

Let C be a closed convex set in \mathbb{R}^n . The distance $f(\mathbf{x}) = \min_{\mathbf{z} \in C} \|\mathbf{z} - \mathbf{x}\|_\dagger$ from \mathbf{x} to C under any norm $\|\mathbf{z}\|_\dagger$ is attained and finite. Furthermore, as observed in Problem 17 of Chap. 6, $f(\mathbf{x})$ is a convex function. In contrast to Euclidean distance, the closest point \mathbf{z} in C to \mathbf{x} may not be unique. Despite

this complication, one can calculate the conjugate $f^*(\mathbf{y})$ and subdifferential $\partial f(\mathbf{x})$. If U is the closed unit ball of the dual norm $\|\mathbf{y}\|_*$, then the conjugate amounts to

$$\begin{aligned}
 f^*(\mathbf{y}) &= \sup_{\mathbf{x}} (\mathbf{y}^* \mathbf{x} - \min_{\mathbf{z} \in C} \|\mathbf{z} - \mathbf{x}\|_{\dagger}) \\
 &= \sup_{\mathbf{x}} \sup_{\mathbf{z} \in C} (\mathbf{y}^* \mathbf{x} - \|\mathbf{z} - \mathbf{x}\|_{\dagger}) \\
 &= \sup_{\mathbf{z} \in C} \{ \mathbf{y}^* \mathbf{z} + \sup_{\mathbf{x}} [\mathbf{y}^* (\mathbf{x} - \mathbf{z}) - \|\mathbf{z} - \mathbf{x}\|_{\dagger}] \} \\
 &= \sup_{\mathbf{z} \in C} [\mathbf{y}^* \mathbf{z} + \sup_{\mathbf{x}} (\mathbf{y}^* \mathbf{x} - \|\mathbf{x}\|_{\dagger})] \\
 &= \sup_{\mathbf{z} \in C} [\mathbf{y}^* \mathbf{z} + \delta_U(\mathbf{y})] \\
 &= \delta_C^*(\mathbf{y}) + \delta_U(\mathbf{y}).
 \end{aligned} \tag{14.7}$$

A subgradient $\mathbf{y} \in \partial f(\mathbf{x})$ is characterized by the equation

$$\mathbf{y}^* \mathbf{x} = f^*(\mathbf{y}) + f(\mathbf{x}) = \delta_C^*(\mathbf{y}) + \delta_U(\mathbf{y}) + f(\mathbf{x}),$$

which clearly forces \mathbf{y} to reside in U . To make further progress, select a point $\mathbf{z} \in C$ with $f(\mathbf{x}) = \|\mathbf{z} - \mathbf{x}\|_{\dagger}$. Because $\delta_C^*(\mathbf{y})$ is the support function of C , this choice yields the inequality

$$\|\mathbf{z} - \mathbf{x}\|_{\dagger} = f(\mathbf{x}) = \mathbf{y}^* \mathbf{x} - \delta_C^*(\mathbf{y}) \leq \mathbf{y}^* (\mathbf{x} - \mathbf{z}). \tag{14.8}$$

Fortunately, the restriction $\mathbf{y} \in U$ and the generalized Cauchy-Schwarz inequality (14.4) imply the opposite inequality, and we conclude that

$$\mathbf{y}^* (\mathbf{x} - \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_{\dagger}. \tag{14.9}$$

The previous example now implies that $\mathbf{y} \in \partial \|\mathbf{x} - \mathbf{z}\|_{\dagger}$. Duplicating the reasoning that led to inequality (14.8) proves that any other point $\mathbf{w} \in C$ satisfies $\mathbf{y}^* (\mathbf{x} - \mathbf{w}) \geq \|\mathbf{x} - \mathbf{z}\|_{\dagger}$. Subtracting this inequality from equality (14.9) gives $\mathbf{y}^* (\mathbf{w} - \mathbf{z}) \leq 0$. Hence, \mathbf{y} belongs to the normal cone $N_C(\mathbf{z})$ to C at \mathbf{z} , and we have proved the containment

$$\partial f(\mathbf{x}) \subset U \cap N_C(\mathbf{z}) \cap \partial \|\mathbf{x} - \mathbf{z}\|_{\dagger}.$$

The reverse containment is also true. Suppose \mathbf{y} belongs to the intersection $U \cap N_C(\mathbf{z}) \cap \partial \|\mathbf{x} - \mathbf{z}\|_{\dagger}$. Then the normal cone condition $\mathbf{y}^* \mathbf{w} \leq \mathbf{y}^* \mathbf{z}$ for all $\mathbf{w} \in C$ implies the equality $\mathbf{y}^* \mathbf{z} = \sup_{\mathbf{w} \in C} \mathbf{y}^* \mathbf{w} = \delta_C^*(\mathbf{y})$. Together with the subgradient condition $\mathbf{y}^* (\mathbf{x} - \mathbf{z}) = \|\mathbf{x} - \mathbf{z}\|_{\dagger}$ of Example 14.4.2, this gives

$$f(\mathbf{x}) = \|\mathbf{z} - \mathbf{x}\|_{\dagger} = \mathbf{y}^* (\mathbf{x} - \mathbf{z}) = \mathbf{y}^* \mathbf{x} - \delta_C^*(\mathbf{y}).$$

In view of the identity $f^*(\mathbf{y}) = \delta_C^*(\mathbf{y}) + \delta_U(\mathbf{y})$ of equation (14.7) and part (b) of Proposition 14.4.4, we are justified in concluding that $\mathbf{y} \in \partial f(\mathbf{x})$.

In the case of the Euclidean norm and a point $\mathbf{x} \notin C$, the point \mathbf{z} is the projection onto C , and the subdifferential $\partial\|\mathbf{x} - \mathbf{z}\|_{\dagger} = \|\mathbf{x} - \mathbf{z}\|^{-1}(\mathbf{x} - \mathbf{z})$ reduces to a unit vector consistent with the normal cone $N_C(\mathbf{z})$. In other words, $f(\mathbf{x})$ is differentiable with gradient $\|\mathbf{x} - \mathbf{z}\|^{-1}(\mathbf{x} - \mathbf{z})$. ■

14.5 The Rules of Convex Differentiation

Some, but not all, of the usual rules of differentiation carry over to convex functions. The simplest rule that successfully generalizes is

$$\partial[\alpha f(\mathbf{x})] = \alpha \partial f(\mathbf{x})$$

for α a positive scalar. This result is just another way of saying that the two supporting plane inequalities

$$\begin{aligned} f(\mathbf{y}) &\geq f(\mathbf{x}) + \mathbf{g}^*(\mathbf{y} - \mathbf{x}) \\ \alpha f(\mathbf{y}) &\geq \alpha f(\mathbf{x}) + (\alpha \mathbf{g})^*(\mathbf{y} - \mathbf{x}) \end{aligned}$$

are logically equivalent. For functions that depend on a single coordinate, another basic rule that the reader can readily verify is

$$\partial f(x_i) = \partial_i f(x_i) \mathbf{e}_i,$$

where ∂ takes the subdifferential with respect to $\mathbf{x} \in \mathbb{R}^n$, ∂_i takes the subdifferential with respect to $x_i \in \mathbb{R}$, and \mathbf{e}_i denotes the i th unit coordinate vector.

Our strategy for deriving other rules is indirect and relies on the characterization (14.5) of the directional derivative as the support function of the subdifferential. Consider a convex function $f(\mathbf{x})$ all of whose directional derivatives $d_{\mathbf{v}}f(\mathbf{z})$ exist and are finite at some point $\mathbf{z} \in \text{dom}(f)$. In this circumstance, equation (14.6) defines the subdifferential $\partial f(\mathbf{z})$ in terms of the directional derivatives $d_{\mathbf{v}}f(\mathbf{z})$. Although definition (14.6) is still operative even when $d_{\mathbf{v}}f(\mathbf{z}) = \infty$ for some \mathbf{v} , we will ignore this possibility. Example 14.3.4 documents the fact that every nonempty closed convex set C is uniquely characterized by its support function

$$\delta_C^*(\mathbf{y}) = \sup_{\mathbf{x} \in C} \mathbf{y}^* \mathbf{x}.$$

We will also need the fact that any nonempty set K with closed convex hull C generates the same support function as C . This is true because

$$\delta_K^*(\mathbf{y}) = \sup_{\mathbf{x} \in K} \mathbf{y}^* \mathbf{x} = \sup_{\mathbf{x} \in C} \mathbf{y}^* \mathbf{x} = \delta_C^*(\mathbf{y}),$$

where the middle equality follows from the identity

$$\max\{a_1, \dots, a_m\} = \max \left\{ \sum_{i=1}^m \lambda_i a_i : \text{all } \lambda_i \geq 0, \sum_{i=1}^m \lambda_i = 1 \right\}$$

and the continuity and linearity of the map $\mathbf{x} \mapsto \mathbf{y}^* \mathbf{x}$.

To generalize the chain rule, consider a convex function $f(\mathbf{y})$ and a differentiable function $g(\mathbf{x})$ whose composition $f \circ g(\mathbf{x})$ with $f(\mathbf{y})$ is convex. If $\mathbf{x} \in \text{dom}(g)$ and $g(\mathbf{x}) \in \text{dom}(f)$, then the difference quotient

$$\begin{aligned} \frac{f \circ g(\mathbf{x} + t\mathbf{v}) - f \circ g(\mathbf{x})}{t} &= \frac{f \circ g(\mathbf{x} + t\mathbf{v}) - f[g(\mathbf{x}) + tdg(\mathbf{x})\mathbf{v}]}{t} \\ &\quad + \frac{f[g(\mathbf{x}) + tdg(\mathbf{x})\mathbf{v}] - f \circ g(\mathbf{x})}{t} \end{aligned}$$

yields in the limit the directional derivative. The first term on the right of this equation vanishes when $g(\mathbf{x})$ is a linear function. The second term tends to $d_{dg(\mathbf{x})\mathbf{v}} f[g(\mathbf{x})]$, the directional derivative of $f(\mathbf{y})$ at the point $\mathbf{y} = g(\mathbf{x})$ in the direction $dg(\mathbf{x})\mathbf{v}$. Even when $g(\mathbf{x})$ is nonlinear, there is still hope that the first term vanishes in the limit. Suppose $g(\mathbf{x})$ belongs to the interior of $\text{dom}(f)$. If we let L be the Lipschitz constant for $f(\mathbf{y})$ near the point $g(\mathbf{x})$ guaranteed by Proposition 6.4.1, then the inequality

$$\begin{aligned} \left| \frac{f \circ g(\mathbf{x} + t\mathbf{v}) - f[g(\mathbf{x}) + tdg(\mathbf{x})\mathbf{v}]}{t} \right| &\leq \frac{L\|g(\mathbf{x} + t\mathbf{v}) - g(\mathbf{x}) - tdg(\mathbf{x})\mathbf{v}\|}{t} \\ &= \frac{o(t)}{t} \end{aligned}$$

shows once again that $d_{\mathbf{v}}(f \circ g)(\mathbf{x}) = d_{dg(\mathbf{x})\mathbf{v}} f(\mathbf{y})$ for $\mathbf{y} = g(\mathbf{x})$. This directional derivative identity translates into the further identity

$$\sup_{\mathbf{z} \in \partial(f \circ g)(\mathbf{x})} \mathbf{v}^* \mathbf{z} = \sup_{\mathbf{z} \in \partial f(\mathbf{y})} \mathbf{v}^* dg(\mathbf{x})^* \mathbf{z} = \sup_{\mathbf{w} \in dg(\mathbf{x})^* \partial f(\mathbf{y})} \mathbf{v}^* \mathbf{w}$$

connecting the corresponding support functions. In view of our earlier remarks regarding support functions and closed convex sets, the convex set $\partial(f \circ g)(\mathbf{x})$ is the closure of the convex set $dg(\mathbf{x})^* \partial f(\mathbf{y})$.

To show that the two sets coincide, it suffices to show that $dg(\mathbf{x})^* \partial f(\mathbf{y})$ is closed. If $\partial f(\mathbf{y})$ is compact, say when \mathbf{y} is an interior point of $\text{dom}(f)$, then $dg(\mathbf{x})^* \partial f(\mathbf{y})$ is compact as well, and the set equality

$$\partial f \circ g(\mathbf{x}) = dg(\mathbf{x})^* \partial f(\mathbf{y})$$

holds. Equality is also true in some circumstances when $\partial f(\mathbf{y})$ is not compact. For example, when $\partial f(\mathbf{y})$ is polyhedral, the image set $\mathbf{A}^* \partial f(\mathbf{y})$ generated by the composition $f(\mathbf{A}\mathbf{x} + \mathbf{b})$ is closed for every compatible matrix \mathbf{A} . A polyhedral set is the nonempty intersection of a finite number of closed halfspaces. Appendix A.3 takes up the subject of polyhedral sets. Proposition A.3.4 is particularly pertinent to the current discussion. The next proposition summarizes our conclusions.

Proposition 14.5.1 (Convex Chain Rule) *Let $f(\mathbf{y})$ be a convex function and $g(\mathbf{x})$ be a differentiable function. If the composition $f \circ g(\mathbf{x})$ is well defined and convex, then the chain rule $\partial f \circ g(\mathbf{x}) = dg(\mathbf{x})^* \partial f(\mathbf{y})$ is*

valid whenever (a) $\mathbf{y} = g(\mathbf{x})$, (b) all possible directional derivatives $d_{\mathbf{v}}f(\mathbf{y})$ exist and are finite, and (c) $dg(\mathbf{x})^*\partial f(\mathbf{y})$ is a closed set.

Proof: See the previous discussion. ■

The composition $f(\mathbf{Ax} + \mathbf{b})$ of $f(\mathbf{y})$ with an affine function $\mathbf{Ax} + \mathbf{b}$ is not the only case of interest. Recall that $f \circ g(\mathbf{x})$ is convex when $f(\mathbf{y})$ is convex and increasing and $g(\mathbf{x})$ is convex. In particular, the composite function $g(\mathbf{x})_+ = \max\{0, g(\mathbf{x})\}$ is convex whenever $g(\mathbf{x})$ is differentiable and convex. Its subdifferential amounts to

$$\partial g(\mathbf{x})_+ = dg(\mathbf{x})^* \begin{cases} 0 & g(\mathbf{x}) < 0 \\ [0, 1] & g(\mathbf{x}) = 0 \\ 1 & g(\mathbf{x}) > 0. \end{cases}$$

In less favorable circumstances, $f \circ g(\mathbf{x})$ is not convex. However, if the directional derivative $d_{\mathbf{v}}(f \circ g)(\mathbf{x})$ is sublinear in \mathbf{v} , then the subdifferential $\partial(f \circ g)(\mathbf{x})$ still makes sense [220].

The sum rule of differentiation is equally easy to prove. Consider two convex functions $f(\mathbf{x})$ and $g(\mathbf{x})$ possessing all possible directional derivatives at the point \mathbf{x} . The identity

$$d_{\mathbf{v}}[f(\mathbf{x}) + g(\mathbf{x})] = d_{\mathbf{v}}f(\mathbf{x}) + d_{\mathbf{v}}g(\mathbf{x})$$

entails the support function identity

$$\sup_{z \in \partial[f(\mathbf{x}) + g(\mathbf{x})]} \mathbf{v}^*z = \sup_{\mathbf{u} \in \partial f(\mathbf{x})} \mathbf{v}^*\mathbf{u} + \sup_{\mathbf{w} \in \partial g(\mathbf{x})} \mathbf{v}^*\mathbf{w} = \sup_{z \in \partial f(\mathbf{x}) + \partial g(\mathbf{x})} \mathbf{v}^*z.$$

It follows that $\partial[f(\mathbf{x}) + g(\mathbf{x})]$ is the closure of the convex set $\partial f(\mathbf{x}) + \partial g(\mathbf{x})$. If either of the two subdifferentials $\partial f(\mathbf{x})$ and $\partial g(\mathbf{x})$ is compact, then the identity $\partial[f(\mathbf{x}) + g(\mathbf{x})] = \partial f(\mathbf{x}) + \partial g(\mathbf{x})$ holds. Indeed, the sum of a closed set and a compact set is always a closed set. Again compactness is not necessary. For example, Proposition A.3.4 demonstrates that the sum of two polyhedral sets is closed. The sum rule is called the Moreau-Rockafellar theorem in the convex calculus literature. Let us again summarize our conclusions by a formal proposition.

Proposition 14.5.2 (Moreau-Rockafellar) *Let $f(\mathbf{x})$ and $g(\mathbf{x})$ be convex functions defined on \mathbb{R}^n . If all possible directional derivatives $d_{\mathbf{v}}f(\mathbf{x})$ and $d_{\mathbf{v}}g(\mathbf{x})$ exist and are finite at a point \mathbf{x} , and the set $\partial f(\mathbf{x}) + \partial g(\mathbf{x})$ is closed, then the sum rule $\partial[f(\mathbf{x}) + g(\mathbf{x})] = \partial f(\mathbf{x}) + \partial g(\mathbf{x})$ is valid.*

Proof: See the foregoing discussion. ■

Example 14.5.1 Mean Value Theorem for Convex Functions.

Let $f(\mathbf{x})$ be a convex function. If two points \mathbf{y} and \mathbf{z} belong to the interior of $\text{dom}(f)$, then the line segment connecting them does as well. Based on

the function $g(t) = f[t\mathbf{y} + (1 - t)\mathbf{z}]$, define the continuous convex function

$$h(t) = g(t) - g(0) - t[g(1) - g(0)].$$

The sum rule and the chain rule imply that $h(t)$ has subdifferential

$$\partial g(t) - g(1) + g(0) = (\mathbf{y} - \mathbf{z})^* \partial f[t\mathbf{y} + (1 - t)\mathbf{z}] - g(1) + g(0).$$

Now $h(0) = h(1) = 0$, so $h(t)$ attains its minimum on the open interval $(0, 1)$. At a minimum point t we have $0 \in \partial h(t)$. It follows that

$$f(\mathbf{y}) - f(\mathbf{z}) = g(1) - g(0) = \mathbf{v}^*(\mathbf{y} - \mathbf{z})$$

for some $\mathbf{v} \in \partial f[t\mathbf{y} + (1 - t)\mathbf{z}]$. ■

Example 14.5.2 *Quantiles and Subdifferentials*

A median μ of n numbers $x_1 \leq x_2 \leq \dots \leq x_n$ satisfies the two inequalities

$$\frac{1}{n} \sum_{x_i \leq \mu} 1 \geq \frac{1}{2} \quad \text{and} \quad \frac{1}{n} \sum_{x_i \geq \mu} 1 \geq \frac{1}{2}.$$

One can relate this to the minimum of the function

$$f(x) = \sum_{i=1}^n |x_i - x|.$$

According to the sum rule, the differential of $f(x)$ equals the set

$$\partial f(x) = - \sum_{x_i > x} 1 + \sum_{x_i = x} [-1, 1] + \sum_{x_i < x} 1.$$

A minimum point μ of $f(x)$ is determined by the condition $0 \in \partial f(\mu)$, which is a disguised form of the median definition just given.

As a generalization take $q \in (0, 1)$ and define

$$\rho_q(r) = \begin{cases} qr & r \geq 0 \\ -(1 - q)r & r < 0. \end{cases}$$

The subdifferential of the function $f_q(x) = \sum_{i=1}^n \rho_q(x_i - x)$ is the set

$$\partial f_q(x) = - \sum_{x_i > x} q + \sum_{x_i = x} [-q, 1 - q] + \sum_{x_i < x} (1 - q).$$

Again the minimum points μ_q of $f_q(x)$ satisfy $0 \in \partial f_q(\mu_q)$, which is equivalent to the q -quantile conditions

$$\frac{1}{n} \sum_{x_i \leq \mu_q} 1 \geq q \quad \text{and} \quad \frac{1}{n} \sum_{x_i \geq \mu_q} 1 \geq 1 - q.$$

Problem 13 of Chap. 8 suggests an MM algorithm for finding μ_q that involves no sorting of the list (x_1, \dots, x_n) . More importantly, one can devise similar MM algorithms for the wider class of quantile regression problems [141]. ■

Example 14.5.3 *Lasso Penalized Estimation*

Lasso penalized estimation minimizes the criterion

$$g(\boldsymbol{\theta}) = f(\boldsymbol{\theta}) + \lambda \sum_{i=1}^p |\theta_i|,$$

where $\lambda \geq 0$ and $f(\boldsymbol{\theta})$ is a convex differentiable loss function. The choice $f(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{y} - \mathbf{X}\boldsymbol{\theta}\|^2$ corresponds to ℓ_2 linear regression. In this case, $\boldsymbol{\theta}$ has $p+1$ components. Penalties are only imposed on the slope parameters $\theta_1, \dots, \theta_p$ and not on the intercept parameter θ_0 . The sum rule implies that $g(\boldsymbol{\theta})$ has subdifferential

$$\partial g(\boldsymbol{\theta}) = \nabla f(\boldsymbol{\theta}) + \lambda \sum_{i=1}^p \partial |\theta_i|.$$

It follows that the stationarity condition $\mathbf{0} \in \partial g(\boldsymbol{\theta})$ amounts to

$$\begin{aligned} \frac{\partial}{\partial \theta_0} f(\boldsymbol{\theta}) &= 0 \\ \frac{\partial}{\partial \theta_i} f(\boldsymbol{\theta}) &\in \lambda \begin{cases} -1 & \theta_i > 0 \\ [-1, 1] & \theta_i = 0 \\ 1 & \theta_i < 0 \end{cases} \end{aligned}$$

for $1 \leq i \leq p$. Hence, a minimum point $\hat{\boldsymbol{\theta}}$ satisfies $|\frac{\partial}{\partial \theta_i} f(\hat{\boldsymbol{\theta}})| \leq \lambda$ for all i . This rather painless deduction extends well beyond ℓ_2 regression. ■

Example 14.5.4 *Euclidean Shrinkage*

Minimization of the strictly convex function

$$f(\mathbf{y}) = \|\mathbf{y}\| + \mathbf{w}^* \mathbf{y} + \frac{\lambda}{2} \|\mathbf{y} - \mathbf{x}\|^2$$

for $\lambda > 0$ illustrates the notion of shrinkage. In the absence of the term $\|\mathbf{y}\|$, the minimum of $f(\mathbf{y})$ occurs at $\mathbf{y} = \lambda^{-1}(\lambda \mathbf{x} - \mathbf{w})$. The presence of the term $\|\mathbf{y}\|$ shrinks the solution to

$$\mathbf{y} = \left[\frac{\|\lambda \mathbf{x} - \mathbf{w}\| - 1}{\|\lambda \mathbf{x} - \mathbf{w}\|} \right]_+ \frac{\lambda \mathbf{x} - \mathbf{w}}{\lambda}.$$

To prove this claim, observe that the subdifferential of $f(\mathbf{y})$ for $\mathbf{y} \neq \mathbf{0}$ collapses to the gradient

$$\nabla f(\mathbf{y}) = \|\mathbf{y}\|^{-1} \mathbf{y} + \mathbf{w} + \lambda(\mathbf{y} - \mathbf{x}).$$

The stationarity condition $\mathbf{0} \in \partial f(\mathbf{y})$ defining a minimum implies

$$(\lambda + \|\mathbf{y}\|^{-1})\mathbf{y} = \lambda\mathbf{x} - \mathbf{w},$$

which in turn implies

$$\|\mathbf{y}\| = \frac{\|\lambda\mathbf{x} - \mathbf{w}\| - 1}{\lambda}.$$

Provided $\|\lambda\mathbf{x} - \mathbf{w}\| > 1$, this is consistent with $\|\mathbf{y}\| > 0$, and we can assert that

$$\mathbf{y} = \|\mathbf{y}\| \frac{\mathbf{y}}{\|\mathbf{y}\|} = \frac{\|\lambda\mathbf{x} - \mathbf{w}\| - 1}{\lambda} \cdot \frac{\lambda\mathbf{x} - \mathbf{w}}{\|\lambda\mathbf{x} - \mathbf{w}\|}$$

furnishes the minimum. If $\|\lambda\mathbf{x} - \mathbf{w}\| \leq 1$, then we must look elsewhere for the minimum. The only other possibility is $\mathbf{y} = \mathbf{0}$. In view of Example 14.4.2, $\partial f(\mathbf{0}) = B + \mathbf{w} - \lambda\mathbf{x}$. The required inclusion $\mathbf{0} \in B + \mathbf{w} - \lambda\mathbf{x}$ now reduces to a tautology. ■

Finally, let us consider a rule with no classical analogue. Generalizing Example 4.4.4, one can easily demonstrate that the maximum function $f(\mathbf{x}) = \max_{1 \leq i \leq p} g_i(\mathbf{x})$ has directional derivative

$$d_{\mathbf{v}}f(\mathbf{x}) = \max_{i \in I(\mathbf{x})} d_{\mathbf{v}}g_i(\mathbf{x}), \tag{14.10}$$

where $I(\mathbf{x})$ is the set of indices i such that $f(\mathbf{x}) = g_i(\mathbf{x})$. All that is required for the validity of formula (14.10) is the upper semicontinuity of the functions $g_i(\mathbf{x})$ at \mathbf{x} and the existence of the directional derivatives $d_{\mathbf{v}}g_i(\mathbf{x})$ for $i \in I(\mathbf{x})$. If we further assume that the $g_i(\mathbf{x})$ are closed and convex, then $f(\mathbf{x})$ is closed and convex, and

$$\sup_{z \in \partial f(\mathbf{x})} \mathbf{v}^*z = \max_{i \in I(\mathbf{x})} \sup_{z_i \in \partial g_i(\mathbf{x})} \mathbf{v}^*z_i.$$

In view of our earlier remarks about support functions and convex hulls, we also have

$$\max_{i \in I(\mathbf{x})} \sup_{z_i \in \partial g_i(\mathbf{x})} \mathbf{v}^*z_i = \sup_{z \in \text{conv}[\cup_{i \in I(\mathbf{x})} \partial g_i(\mathbf{x})]} \mathbf{v}^*z.$$

If each subdifferential $\partial g_i(\mathbf{x})$ is compact, then the finite union $\cup_{i \in I(\mathbf{x})} \partial g_i(\mathbf{x})$ is also compact. Because the convex hull of a compact set is compact (Proposition 6.2.4), the conclusion

$$\partial f(\mathbf{x}) = \text{conv}[\cup_{i \in I(\mathbf{x})} \partial g_i(\mathbf{x})] \tag{14.11}$$

emerges. This formula is valid in more general circumstances. For example, if i is a continuous variable, then it suffices for the $g_i(\mathbf{x})$ to be convex, the set of indices i to be compact, and the function $i \mapsto g_i(\mathbf{x})$ to be upper semicontinuous for each \mathbf{x} [131, 226]. Proposition A.6.6 in the Appendix also addresses this topic.

Example 14.5.5 *Minima of Max Functions*

When the functions $g_i(\mathbf{x})$ are convex, a point \mathbf{y} minimizes

$$f(\mathbf{x}) = \max_{1 \leq i \leq n} g_i(\mathbf{x})$$

if and only if $\mathbf{0} \in \partial f(\mathbf{y})$. Equivalently, $\mathbf{0}$ belongs to the convex hull of the vectors $\{\nabla g_i(\mathbf{y})\}_{i \in I(\mathbf{y})}$. According to Proposition 5.3.2, a necessary and sufficient condition for the latter event to occur is that no vector \mathbf{u} exists with $dg_i(\mathbf{y})\mathbf{u} < 0$ for all $i \in I(\mathbf{y})$. Of course, such a vector \mathbf{u} would constitute a descent direction along which $f(\mathbf{x})$ could be locally reduced. One can reformulate the problem of minimization of $f(\mathbf{x})$ as a linear program when the $g_i(\mathbf{x}) = \mathbf{w}_i^* \mathbf{x} + a_i$ are affine. In this case one just minimizes the scalar t subject to the inequality constraints $\mathbf{w}_i^* \mathbf{x} + a_i \leq t$ for all i . ■

Example 14.5.6 *Sum of the m Largest Order Statistics*

Consider the order statistics $x_{(1)} \leq \dots \leq x_{(n)}$ corresponding to a point $\mathbf{x} \in \mathbb{R}^n$. The sum

$$s_m(\mathbf{x}) = \sum_{i=n-m+1}^n x_{(i)} = \max_{|T|=m} \sum_{i \in T} x_i$$

of the m largest order statistics is a convex function of \mathbf{x} . Here T is any subset of $\{1, \dots, n\}$ of size $|T| = m$. The gradient of $\sum_{i \in T} x_i$ is the vector $\mathbf{1}_T$ whose i th entry is 1 if $i \in T$ and 0 otherwise. The subdifferential $\partial s_m(\mathbf{x})$ is therefore

$$\text{conv} \{ \mathbf{1}_T : \sum_{i \in T} x_i = s_m(\mathbf{x}), |T| = m \}$$

When all of the components of \mathbf{x} are unique, $s_m(\mathbf{x})$ is differentiable. ■

Example 14.5.7 *Maximum Eigenvalue of a Symmetric Matrix*

As mentioned in Example 6.3.8, the maximum eigenvalue $\lambda_{\max}(\mathbf{M})$ of a symmetric matrix is a convex function whose value is determined by the maximum Rayleigh quotient

$$\lambda_{\max}(\mathbf{M}) = \sup_{\|\mathbf{x}\|=1} \mathbf{x}^* \mathbf{M} \mathbf{x} = \sup_{\|\mathbf{x}\|=1} \text{tr}(\mathbf{M} \mathbf{x} \mathbf{x}^*).$$

The Frobenius inner product map $\mathbf{M} \mapsto \text{tr}(\mathbf{M} \mathbf{x} \mathbf{x}^*)$ is linear in \mathbf{M} and has differential $\mathbf{x} \mathbf{x}^*$. Hence, the subdifferential $\partial \lambda_{\max}(\mathbf{M})$ is the convex hull of the set $\{ \mathbf{x} \mathbf{x}^* : \|\mathbf{x}\| = 1, \mathbf{M} \mathbf{x} = \lambda_{\max}(\mathbf{M}) \mathbf{x} \}$. When there is a unique unit eigenvector \mathbf{x} up to sign, $\lambda_{\max}(\mathbf{M})$ is differentiable with gradient $\mathbf{x} \mathbf{x}^*$. ■

14.6 Spectral Functions

We are already acquainted with some of the connections between eigenvalues and convexity. The lovely theory of Lewis [17] extends these previously scattered results. Once again the Fenchel conjugate is the tool of choice. Lewis' theory covers the composition of a symmetric function $f : \mathbb{R}^n \mapsto \mathbb{R}$ with the eigenvalue map $\lambda(\mathbf{X})$ defined on $n \times n$ symmetric matrices \mathbf{X} . Recall that $f(\mathbf{x})$ is symmetric if $f(\mathbf{w}) = f(\mathbf{x})$ whenever the vector arguments \mathbf{w} and \mathbf{x} agree up to a permutation of their entries. The vector $\lambda(\mathbf{X})$ presents the spectrum of \mathbf{X} ordered from the largest to the smallest eigenvalue. The next proposition shows how to calculate the Fenchel conjugate of the composite function $f \circ \lambda(\mathbf{X})$. In the proposition the operator $\text{diag}(\mathbf{x})$ promotes a vector \mathbf{x} to a diagonal matrix with x_i as its i th diagonal entry.

Proposition 14.6.1 *Suppose $f(\mathbf{x})$ has Fenchel conjugate $f^*(\mathbf{y})$. Then the spectral function $g(\mathbf{X}) = f \circ \lambda(\mathbf{X})$ has Fenchel conjugate $f^* \circ \lambda(\mathbf{Y})$. Hence, if $f(\mathbf{x})$ is closed and convex, then $f \circ \lambda(\mathbf{X})$ is closed and convex and equals its Fenchel biconjugate.*

Proof: Verification relies heavily on Fan's inequality (A.4.2) proved in Appendix A.4. If we denote the set of $n \times n$ symmetric matrices by S^n , then Fan's inequality implies

$$\begin{aligned} (f \circ \lambda)^*(\mathbf{Y}) &= \sup_{\mathbf{X} \in S^n} [\text{tr}(\mathbf{Y}\mathbf{X}) - f \circ \lambda(\mathbf{X})] \\ &\leq \sup_{\mathbf{X} \in S^n} [\lambda(\mathbf{Y})^* \lambda(\mathbf{X}) - f \circ \lambda(\mathbf{X})] \\ &\leq \sup_{\mathbf{x} \in \mathbb{R}^n} [\lambda(\mathbf{Y})^* \mathbf{x} - f(\mathbf{x})] \\ &= f^*[\lambda(\mathbf{Y})]. \end{aligned}$$

The reverse inequality follows from the spectral decomposition

$$\mathbf{Y} = \mathbf{U}^* \text{diag}(\mathbf{y})\mathbf{U}$$

of \mathbf{Y} in terms of an orthogonal matrix \mathbf{U} . This leads to

$$\begin{aligned} f^*[\lambda(\mathbf{Y})] &= \sup_{\mathbf{x} \in \mathbb{R}^n} [\lambda(\mathbf{Y})^* \mathbf{x} - f(\mathbf{x})] \\ &= \sup_{\mathbf{x} \in \mathbb{R}^n} \{ \text{tr}[\mathbf{U}\mathbf{Y}\mathbf{U}^* \text{diag}(\mathbf{x})] - f(\mathbf{x}) \} \\ &= \sup_{\mathbf{x} \in \mathbb{R}^n} \{ \text{tr}[\mathbf{Y}\mathbf{U}^* \text{diag}(\mathbf{x})\mathbf{U}] - f[\lambda(\mathbf{U}^* \text{diag}(\mathbf{x})\mathbf{U})] \} \\ &\leq \sup_{\mathbf{X} \in S^n} \{ \text{tr}(\mathbf{Y}^* \mathbf{X}) - f[\lambda(\mathbf{X})] \} \\ &= (f \circ \lambda)^*(\mathbf{Y}). \end{aligned}$$

Thus, the two extremes of each inequality are equal. The second claim of the proposition follows from a double application of Proposition 14.3.2. ■

The maximum eigenvalue $\lambda(\mathbf{M})$ of a symmetric matrix is a spectral function. We have already calculated its subdifferential. The next proposition allows us to calculate the subdifferentials of more complicated spectral functions.

Proposition 14.6.2 *Suppose $f(\mathbf{x})$ is closed and convex. A symmetric matrix \mathbf{Y} belongs to the subdifferential $\partial(f \circ \lambda)(\mathbf{X})$ if and only if the ordered vector $\lambda(\mathbf{Y})$ belongs to the subdifferential $\partial f[\lambda(\mathbf{X})]$ and \mathbf{X} and \mathbf{Y} have simultaneous ordered spectral decompositions $\mathbf{X} = \mathbf{U}^* \text{diag} \circ \lambda(\mathbf{X}) \mathbf{U}$ and $\mathbf{Y} = \mathbf{U}^* \text{diag} \circ \lambda(\mathbf{Y}) \mathbf{U}$.*

Proof: The combination of Proposition 14.6.1, the Fenchel-Young inequality, and Fan’s inequality yields

$$\begin{aligned} (f \circ \lambda)(\mathbf{X}) + (f \circ \lambda)^*(\mathbf{Y}) &= f[\lambda(\mathbf{X})] + f^*[\lambda(\mathbf{Y})] \\ &\geq \lambda(\mathbf{X})^* \lambda(\mathbf{Y}) \\ &\geq \text{tr}(\mathbf{YX}). \end{aligned}$$

The matrix \mathbf{Y} belongs to $\partial(f \circ \lambda)(\mathbf{X})$ if and only if the extreme members of these two inequalities agree. The first inequality is an equality if and only if $\lambda(\mathbf{Y})$ belongs to the subdifferential $\partial f[\lambda(\mathbf{X})]$. As observed in Proposition A.4.2, the second inequality is an equality if and only if \mathbf{X} and \mathbf{Y} have simultaneous ordered spectral decompositions. ■

Example 14.6.1 *Sample Calculations with Spectral Functions*

Various matrix norms can be identified as spectral functions. For instance, the nuclear and Euclidean matrix norms originate from the symmetric convex functions

$$f_1(\mathbf{x}) = \sum_{i=1}^n |x_i| \quad \text{and} \quad f_2(\mathbf{x}) = \max_{1 \leq i \leq n} |x_i|.$$

Example 14.5.7 and Proposition 14.6.2 clearly produce the same subdifferential for the Euclidean matrix norm. The spectral function

$$f_3(\mathbf{x}) = \sum_{i=1}^n \max\{-x_i, 0\},$$

arises when negative eigenvalues are downplayed.

The forgoing theory can be applied to penalized approximation of symmetric matrices. Indeed, consider minimization of

$$\frac{1}{2} \|\mathbf{Y} - \mathbf{X}\|_F^2 + \rho f_i \circ \lambda(\mathbf{Y}) \tag{14.12}$$

as a function of the symmetric matrix \mathbf{Y} for one of the three functions just defined. Fermat’s rule requires

$$\mathbf{0} \in \mathbf{Y} - \mathbf{X} + \rho \partial(f_i \circ \lambda)(\mathbf{Y}).$$

In view of Proposition 14.6.2,

$$\mathbf{X} \in \mathbf{U}^* \text{diag}(\mathbf{y})\mathbf{U} + \rho\mathbf{U}^* \text{diag}[\partial f_i(\mathbf{y})]\mathbf{U}$$

for $\mathbf{y} = \lambda(\mathbf{Y})$ and some orthogonal matrix \mathbf{U} . It follows that \mathbf{X} and \mathbf{Y} have simultaneous spectral decompositions under \mathbf{U} . If we let $\mathbf{x} = \lambda(\mathbf{X})$ and multiply the relation

$$\mathbf{0} \in \mathbf{U}^* \text{diag}(\mathbf{y})\mathbf{U} - \mathbf{U}^* \text{diag}(\mathbf{x})\mathbf{U} + \rho\mathbf{U}^* \text{diag}[\partial f_i(\mathbf{y})]\mathbf{U},$$

by \mathbf{U} on the left and \mathbf{U}^* on the right, then the relation

$$\mathbf{0} \in \mathbf{y} - \mathbf{x} + \rho\partial f_i(\mathbf{y}),$$

emerges. But this is just Fermat's rule for minimizing the spectral function $g(\mathbf{y}) = \frac{1}{2}\|\mathbf{y} - \mathbf{x}\|^2 + \rho f_i(\mathbf{y})$.

The problem of minimizing $g(\mathbf{y})$ is separable for the penalty $f_1(\mathbf{y})$. Its solution has components

$$y_i = \begin{cases} x_i - \rho & x_i > \rho \\ 0 & |x_i| \leq \rho \\ x_i + \rho & x_i < -\rho \end{cases} \quad (14.13)$$

exhibiting shrinkage. For the penalty $f_3(\mathbf{y})$, the minimum point of $g(\mathbf{y})$ has a solution with components

$$y_i = \begin{cases} x_i & x_i > 0 \\ 0 & -\rho \leq x_i \leq 0 \\ x_i + \rho & x_i < -\rho \end{cases}$$

exhibiting one-sided shrinkage.

Finding the minimum of $g(\mathbf{y})$ for the $f_2(\mathbf{y})$ penalty is harder because the objective function is no longer separable. Inspection of $g(\mathbf{y})$ shows that the solution must satisfy $\text{sgn}(y_i) = \text{sgn}(x_i)$ and $|y_i| \leq |x_i|$ for all i . Thus, there is no loss in generality in assuming that $0 \leq y_i \leq x_i$ for all i . Instead of exploiting subdifferentials directly, let us focus on forward directional derivatives. At the point $\mathbf{y} = \mathbf{0}$ an easy calculation shows that

$$d_{\mathbf{v}}g(\mathbf{0}) = \sum_{i=1}^n (0 - x_i)v_i + \rho \max_{1 \leq i \leq n} |v_i|.$$

According to Proposition 6.5.2, $\mathbf{0}$ furnishes the minimum of $g(\mathbf{y})$ if and only if all of these directional derivatives are nonnegative. In view of the generalized Cauchy-Schwarz inequality

$$\left| \sum_{i=1}^n (0 - x_i)v_i \right| \leq \|\mathbf{x}\|_1 \|\mathbf{v}\|_\infty,$$

the vector $\mathbf{0}$ qualifies provided $\|\mathbf{x}\|_1 \leq \rho$.

When $\|\mathbf{x}\|_1 > \rho$, there is a simple recipe for constructing the solution \mathbf{y} . Gradually decrease $r = \|\mathbf{y}\|_\infty$ from $\|\mathbf{x}\|_\infty$ to 0 until the condition

$$\sum_{x_i \geq r} (x_i - r) = \rho$$

is met. We claim that the vector \mathbf{y} with components

$$y_i = \begin{cases} r & x_i \geq r \\ x_i & 0 \leq x_i < r \end{cases}$$

provides the minimum. It suffices to prove that the directional derivative

$$d_{\mathbf{v}}g(\mathbf{y}) = \sum_{x_i \geq r} (y_i - x_i)v_i + \rho \max_{x_i \geq r} v_i$$

is nonnegative regardless of how we choose \mathbf{v} . But this fact follows from the inequality

$$\begin{aligned} \sum_{x_i \geq r} (y_i - x_i)v_i &= \sum_{x_i \geq r} (x_i - y_i)(-v_i) \\ &\geq \min_{x_i \geq r} (-v_i) \sum_{x_i \geq r} (x_i - y_i) \\ &= \rho \min_{x_i \geq r} (-v_i) \\ &= -\rho \max_{x_i \geq r} v_i. \end{aligned}$$

Given this solution for \mathbf{x} with nonnegative components, it is straightforward to recover the general solution. \blacksquare

Example 14.6.2 Matrix Completion Problem

Similar calculations are pertinent to matrix completion [37]. Consider a matrix $\mathbf{Y} = (y_{ij})$, some of whose entries are unobserved. If Δ denotes the set of index pairs (i, j) such that y_{ij} is observed, then it is convenient to define the projected matrix $P_\Delta(\mathbf{Y})$ with entries

$$P_\Delta(\mathbf{Y}) = \begin{cases} y_{ij} & (i, j) \in \Delta \\ 0 & (i, j) \notin \Delta. \end{cases}$$

Its orthogonal complement $P_\Delta^\perp(\mathbf{Y})$ satisfies $P_\Delta^\perp(\mathbf{Y}) + P_\Delta(\mathbf{Y}) = \mathbf{Y}$. In the matrix completion problem one seeks to minimize the criterion

$$\frac{1}{2} \sum_{(i,j) \in \Delta} (y_{ij} - x_{ij})^2 + \rho \|\mathbf{X}\|_*$$

involving the nuclear norm of $\mathbf{X} = (x_{ij})$. One way of attacking the problem is to majorize the objective function at the current iterate \mathbf{X}_m by

$$g(\mathbf{X} \mid \mathbf{X}_m) = \frac{1}{2} \|P_\Delta(\mathbf{Y}) + P_\Delta^\perp(\mathbf{X}_m) - \mathbf{X}\|_F^2 + \rho \|\mathbf{X}\|_*.$$

In light of the MM principle, minimizing $g(\mathbf{X} \mid \mathbf{X}_m)$ drives the matrix completion criterion downhill [189].

Now set $\mathbf{Z}_m = P_\Delta(\mathbf{Y}) + P_\Delta^\perp(\mathbf{X}_m)$ and contemplate the expansion

$$g(\mathbf{X} \mid \mathbf{X}_m) = \frac{1}{2} \|\mathbf{Z}_m\|_F^2 - \text{tr}(\mathbf{Z}_m^* \mathbf{X}) + \frac{1}{2} \|\mathbf{X}\|_F^2 + \rho \|\mathbf{X}\|_*.$$

Suppose \mathbf{X} has svd $\mathbf{P}\Sigma\mathbf{Q}^*$ and \mathbf{Z}_m has svd $\mathbf{U}\Omega\mathbf{V}^*$. According to Fan's inequality as discussed in Appendix A.5, $-\text{tr}(\mathbf{Z}_m^* \mathbf{X}) \geq -\sum_i \omega_i \sigma_i$, with equality when $\mathbf{U} = \mathbf{P}$ and $\mathbf{V} = \mathbf{Q}$. Here, the ω_i and σ_i are the ordered singular values of \mathbf{X} and \mathbf{Z}_m , respectively. Furthermore, neither the Frobenius norm $\|\mathbf{X}\|_F = (\sum_i \sigma_i^2)^{1/2}$ nor the nuclear norm $\|\mathbf{X}\|_* = \sum_i \sigma_i$ depends on the orthogonal matrices \mathbf{P} and \mathbf{Q} of singular vectors. Hence, it is optimal to take $\mathbf{P} = \mathbf{U}$ and $\mathbf{Q} = \mathbf{V}$. The problem therefore becomes one of minimizing

$$\frac{1}{2} \sum_i (\omega_i - \sigma_i)^2 + \rho \sum_i \sigma_i$$

subject to the nonnegativity constraints $\sigma_i \geq 0$. The solution is given by equation (14.13) with $x_i = \omega_i \geq 0$ and $y_i = \sigma_i$. In practice only the largest singular values of Ω need be extracted. The Lanczos procedure [107] efficiently computes these and their corresponding singular vectors.

A related problem is to minimize $\|\mathbf{X}\|_*$ subject to the quadratic constraint $\frac{1}{2} \sum_{(i,j) \in \Delta} (y_{ij} - x_{ij})^2 \leq \epsilon$ [35]. This problem can be recast as minimizing the penalized convex function $\|\mathbf{X}\|_* + \delta_C(\mathbf{X})$, where C is the closed convex set

$$C = \left\{ \mathbf{X} = (x_{ij}) : \frac{1}{2} \sum_{(i,j) \in \Delta} (y_{ij} - x_{ij})^2 \leq \epsilon \right\}.$$

To derive an MM algorithm, suppose that the current iterate is \mathbf{X}_m and that $\mathbf{Z}_m = P_\Delta(\mathbf{Y}) + P_\Delta^\perp(\mathbf{X}_m)$. If we define the closed convex set

$$C_m = \left\{ \mathbf{X} = (x_{ij}) : \frac{1}{2} \sum_i \sum_j (z_{mij} - x_{ij})^2 \leq \epsilon \right\},$$

then it is obvious that $\|\mathbf{X}\|_* + \delta_{C_m}(\mathbf{X})$ majorizes $\|\mathbf{X}\|_* + \delta_C(\mathbf{X})$ around the point \mathbf{X}_m . In this majorization we allow infinite function values.

Minimization of the surrogate function $\|\mathbf{X}\|_* + \delta_{C_m}(\mathbf{X})$ again relies on the fact that the nuclear and Frobenius norms are invariant under left and right multiplication of their arguments by orthogonal matrices. As we have just argued, we can assume that \mathbf{X} and \mathbf{Z}_m have svds $\mathbf{U}\Sigma\mathbf{V}^*$ and $\mathbf{U}\Omega\mathbf{V}^*$ involving shared orthogonal matrices \mathbf{U} and \mathbf{V} . Thus, the current problem reduces to minimizing $\sum_i \sigma_i$ subject to $\frac{1}{2} \sum_i (\omega_i - \sigma_i)^2 \leq \epsilon$ and $\sigma_i \geq 0$ for all i . Of course, the singular values ω_i are nonnegative as well. If $\frac{1}{2} \sum_i \omega_i^2 \leq \epsilon$, then the trivial solution $\sigma_i = 0$ for all i holds.

Hence, suppose that $\frac{1}{2} \sum_i \omega_i^2 > \epsilon$ and form the Lagrangian

$$\mathcal{L}(\boldsymbol{\sigma}, \lambda, \boldsymbol{\mu}) = \sum_i \sigma_i + \lambda \left[\frac{1}{2} \sum_i (\omega_i - \sigma_i)^2 - \epsilon \right] - \sum_i \mu_i \sigma_i.$$

If we assume $\frac{1}{2} \sum_i (\omega_i - \sigma_i)^2 = \epsilon$ and $\lambda > 0$, then the stationarity condition

$$0 = 1 + \lambda(\sigma_i - \omega_i) - \mu_i$$

and complementary slackness imply

$$\sigma_i = \begin{cases} \omega_i - \lambda^{-1} & \omega_i - \lambda^{-1} > 0 \\ 0 & \omega_i - \lambda^{-1} \leq 0. \end{cases}$$

The condition $\frac{1}{2} \sum_i (\omega_i - \sigma_i)^2 = \epsilon$ then amounts to the identity

$$\frac{1}{2} \sum_i \min\{\omega_i, \lambda^{-1}\}^2 = \epsilon.$$

The continuous function $f(\beta) = \frac{1}{2} \sum_i \min\{\omega_i, \beta\}^2$ is strictly increasing on the interval $[0, \max_i \omega_i]$. Since $f(0) = 0$ and $f(\max_i \omega_i) > \epsilon$ by assumption, there is a unique positive solution. For n singular values, the solution belongs to the interval $[\omega_k, \omega_{k+1}]$ if and only if

$$\sum_{i=1}^k \omega_i^2 + (n-k)\omega_k^2 \leq 2\epsilon \leq \sum_{i=1}^k \omega_i^2 + (n-k)\omega_{k+1}^2.$$

Because this is equivalent to

$$(n-k)\omega_k^2 \leq 2\epsilon - \sum_i \omega_i^2 + \sum_{i>k} \omega_i^2 \leq (n-k)\omega_{k+1}^2$$

and $\sum_i \omega_i^2 = \|\mathbf{Z}_m\|_F^2$, the solution again depends on only the largest singular values. \blacksquare

Example 14.6.3 *Stable Estimation of a Covariance Matrix*

Example 6.5.7 demonstrates that the sample mean and sample covariance matrix

$$\begin{aligned} \bar{\mathbf{y}} &= \frac{1}{k} \sum_{j=1}^k \mathbf{y}_j \\ \mathbf{S} &= \frac{1}{k} \sum_{j=1}^k (\mathbf{y}_j - \bar{\mathbf{y}})(\mathbf{y}_j - \bar{\mathbf{y}})^* \end{aligned}$$

are the maximum likelihood estimates of the theoretical mean $\boldsymbol{\mu}$ and theoretical covariance $\boldsymbol{\Omega}$ of a random sample $\mathbf{y}_1, \dots, \mathbf{y}_k$ from a multivariate

normal distribution. When the number of components n of \mathbf{y} exceeds the sample size k , this analysis breaks down because it assumes that \mathbf{S} is invertible. To deal with this dilemma and to stabilize estimation generally, we add a penalty. The estimate of $\boldsymbol{\mu}$ remains $\bar{\mathbf{y}}$.

As in Sect. 13.4, we impose a prior $p(\boldsymbol{\Omega})$ on $\boldsymbol{\Omega}$. Now, however, the prior is designed to steer the eigenvalues of $\boldsymbol{\Omega}$ away from the extremes of 0 and ∞ . The reasonable choice

$$p(\boldsymbol{\Omega}) \propto e^{-\frac{\lambda}{2}[\alpha\|\boldsymbol{\Omega}\|_* + (1-\alpha)\|\boldsymbol{\Omega}^{-1}\|_*]},$$

relies on the nuclear norms of $\boldsymbol{\Omega}$ and $\boldsymbol{\Omega}^{-1}$, a positive strength constant λ , and an admixture constant $\alpha \in (0, 1)$. This is a proper prior on the set of invertible matrices because

$$e^{-\frac{\lambda}{2}[\alpha\|\boldsymbol{\Omega}\|_* + (1-\alpha)\|\boldsymbol{\Omega}^{-1}\|_*]} \leq e^{-c\lambda\|\boldsymbol{\Omega}\|_F}$$

for some positive constant c by virtue of the equivalence of any two vector norms on \mathbf{R}^{n^2} . The normalizing constant of $p(\boldsymbol{\Omega})$ is irrelevant in the ensuing discussion. Consider therefore minimization of the function

$$f(\boldsymbol{\Omega}) = \frac{k}{2} \ln \det \boldsymbol{\Omega} + \frac{k}{2} \operatorname{tr}(\mathbf{S}\boldsymbol{\Omega}^{-1}) + \frac{\lambda}{2} [\alpha\|\boldsymbol{\Omega}\|_* + (1-\alpha)\|\boldsymbol{\Omega}^{-1}\|_*].$$

The maximum of $-f(\boldsymbol{\Omega})$ occurs at the posterior mode. In the limit as λ tends to 0, $-f(\boldsymbol{\Omega})$ reduces to the loglikelihood.

Fortunately, three of the four terms of $f(\boldsymbol{\Omega})$ can be expressed as functions of the eigenvalues e_i of $\boldsymbol{\Omega}$. The trace contribution presents a greater challenge. Let $\mathbf{S} = \mathbf{U}\mathbf{D}\mathbf{U}^*$ denote the spectral decomposition of \mathbf{S} with nonnegative diagonal entries d_i ordered from largest to smallest. Likewise, let $\boldsymbol{\Omega} = \mathbf{V}\mathbf{E}\mathbf{V}^*$ denote the spectral decomposition of $\boldsymbol{\Omega}$ with positive diagonal entries e_i ordered from largest to smallest. In view of Fan's inequality (A.6), we can assert that

$$-\operatorname{tr}(\mathbf{S}\boldsymbol{\Omega}^{-1}) \leq -\sum_{i=1}^n \frac{d_i}{e_i},$$

with equality if and only if $\mathbf{V} = \mathbf{U}$. Consequently, we make the latter assumption and replace $f(\boldsymbol{\Omega})$ by

$$g(\mathbf{E}) = \frac{k}{2} \sum_{i=1}^n \ln e_i + \frac{k}{2} \sum_{i=1}^n \frac{d_i}{e_i} + \frac{\lambda}{2} \left[\alpha \sum_{i=1}^n e_i + (1-\alpha) \sum_{i=1}^n \frac{1}{e_i} \right].$$

At a stationary point of $g(\mathbf{E})$, we have

$$0 = \frac{k}{e_i} - \frac{k d_i + \lambda(1-\alpha)}{e_i^2} + \lambda \alpha.$$

The solution to this essentially quadratic equation is

$$e_i = \frac{-k + \sqrt{k^2 + 4\lambda\alpha[kd_i + \lambda(1 - \alpha)]}}{2\lambda\alpha}. \quad (14.14)$$

We reject the negative root as inconsistent with $\mathbf{\Omega}$ being positive definite. For the special case $k = 0$ of no data, all $e_i = \sqrt{(1 - \alpha)/\alpha}$, and the prior mode occurs at a multiple of the identity matrix.

Holding all but one variable fixed in formula (14.14), one can demonstrate after a fair amount of algebra that

$$e_i = d_i + \frac{\lambda(1 - \alpha - \alpha d_i^2)}{k} + O\left(\frac{1}{k^2}\right), \quad k \rightarrow \infty$$

$$e_i = \sqrt{\frac{1 - \alpha}{\alpha}} + \left[\sqrt{\frac{1 - \alpha}{\alpha}} \frac{kd_i}{2(1 - \alpha)} - \frac{k}{2\alpha} \right] \frac{1}{\lambda} + O\left(\frac{1}{\lambda^2}\right), \quad \lambda \rightarrow \infty.$$

These asymptotic expansions accord with common sense. Namely, the data eventually overwhelms a fixed prior, and increasing the penalty strength for a fixed amount of data pulls the estimate of $\mathbf{\Omega}$ toward the prior mode. Choice of the constants λ and α is an issue. To match the prior to the scale of the data, we recommend determining α as the solution to the equation

$$n\sqrt{\frac{1 - \alpha}{\alpha}} = \text{tr} \left(\sqrt{\frac{1 - \alpha}{\alpha}} \mathbf{I} \right) = \text{tr}(\mathbf{S}).$$

Cross validation leads to a reasonable choice of λ . For the sake of brevity, we omit further details. For a summary of other approaches to this subject, consult the reference [173]. ■

14.7 A Convex Lagrange Multiplier Rule

The Lagrange multiplier rule is one of the dominant themes of optimization theory. In convex programming it represents both a necessary and sufficient condition for a minimum point. Our previous proofs of the rule invoked differentiability. In fact, differentiability assumptions can be dismissed in deriving the multiplier rule for convex programs. Recall that we posed convex programming as the problem of minimizing a convex function $f(\mathbf{y})$ subject to the constraints

$$g_i(\mathbf{y}) = 0, \quad 1 \leq i \leq p, \quad h_j(\mathbf{y}) \leq 0, \quad 1 \leq j \leq q,$$

where the $g_i(\mathbf{y})$ are affine functions and the $h_j(\mathbf{y})$ are convex functions. One can simplify the statement of many convex programs by intersecting the essential domains of the objective and constraint functions with a

closed convex set C . For example, although the set of positive semidefinite matrices is convex, it is awkward to represent it as an intersection of convex sets determined by affine equality constraints and simple convex inequality constraints.

In proving the multiplier rule anew, we will call on Slater’s constraint qualification. In the current context, this entails postulating the existence of a point $\mathbf{z} \in C$ such that $g_i(\mathbf{z}) = 0$ for all i and $h_j(\mathbf{z}) < 0$ for all j . In addition we assume that the constraint gradient vectors $\nabla g_i(\mathbf{y})$ are linearly independent. These preliminaries put us into position to restate and prove the Lagrange multiplier rule for convex programs.

Proposition 14.7.1 *Suppose that $f(\mathbf{y})$ achieves its minimum subject to the constraints at the point $\mathbf{x} \in C$. Then there exists a Lagrangian function*

$$\mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \lambda_0 f(\mathbf{y}) + \sum_{i=1}^p \lambda_i g_i(\mathbf{y}) + \sum_{j=1}^q \mu_j h_j(\mathbf{y})$$

characterized by the following three properties: (a) $\mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ achieves its minimum over C at \mathbf{x} , (b) the multipliers λ_0 and μ_j are nonnegative, and (c) the complementary slackness conditions $\mu_j h_j(\mathbf{x}) = 0$ hold. If Slater’s constraint qualification is true and \mathbf{x} is an interior point of C , then one can take $\lambda_0 = 1$. Conversely, if properties (a) through (c) hold with $\lambda_0 = 1$, then \mathbf{x} minimizes $f(\mathbf{y})$ subject to the constraints.

Proof: To prove the necessity of the three properties, we separate a convex set and a point by a hyperplane. Accordingly, define the set S to consist of all points $\mathbf{u} \in \mathbb{R}^{p+q+1}$ such that for some $\mathbf{y} \in C$, we have $u_0 \geq f(\mathbf{y})$, $u_i = g_i(\mathbf{y})$ for $1 \leq i \leq p$, and $u_{p+j} \geq h_j(\mathbf{y})$ for $1 \leq j \leq q$. To show that the set S is convex, suppose $\mathbf{y} \in C$ corresponds to $\mathbf{u} \in S$, $\mathbf{z} \in C$ corresponds to $\mathbf{v} \in S$, and α and β are nonnegative numbers summing to 1. The relations

$$\begin{aligned} f(\alpha\mathbf{y} + \beta\mathbf{z}) &\leq \alpha f(\mathbf{y}) + \beta f(\mathbf{z}) \leq \alpha u_0 + \beta v_0 \\ g_i(\alpha\mathbf{y} + \beta\mathbf{z}) &= \alpha g_i(\mathbf{y}) + \beta g_i(\mathbf{z}) = \alpha u_i + \beta v_i, \quad 1 \leq i \leq p \\ h_j(\alpha\mathbf{y} + \beta\mathbf{z}) &\leq \alpha h_j(\mathbf{y}) + \beta h_j(\mathbf{z}) \leq \alpha u_{p+j} + \beta v_{p+j}, \quad 1 \leq j \leq q \end{aligned}$$

and the convexity of C prove the convexity claim. The point to be separated from S is $[f(\mathbf{x}), \mathbf{0}^*]^*$. It belongs to S because \mathbf{x} is feasible. It lies on the boundary of S because the point $[f(\mathbf{x}) - \epsilon, \mathbf{0}^*]^*$ does not belong to S for any $\epsilon > 0$.

Application of Proposition 6.2.3 shows that there exists a nontrivial vector $\boldsymbol{\omega}$ such that $\omega_0 f(\mathbf{x}) \leq \boldsymbol{\omega}^* \mathbf{u}$ for all $\mathbf{u} \in S$. Identify the entries of $\boldsymbol{\omega}$ with $\lambda_0, \lambda_1, \dots, \lambda_p, \mu_1, \dots, \mu_q$, in that order. Sending u_0 to ∞ implies $\lambda_0 \geq 0$; similarly, sending u_{p+j} to ∞ implies $\mu_j \geq 0$. If $h_j(\mathbf{x}) < 0$, then the vector $\mathbf{u} \in S$ with $f(\mathbf{x})$ as entry 0, $h_j(\mathbf{x})$ as entry $p+j$, and 0’s elsewhere demonstrates that $\mu_j = 0$. This proves properties (b) and (c). To verify property

(a), take $\mathbf{y} \in C$ and put $u_0 = f(\mathbf{y})$, $u_i = g_i(\mathbf{y})$, and $u_{p+j} = h_j(\mathbf{y})$. Then $\mathbf{u} \in S$, and the separating hyperplane condition reads

$$\mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \geq \lambda_0 f(\mathbf{x}) = \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}),$$

proving property (a).

Next suppose $\lambda_0 = 0$ and \mathbf{z} is a Slater point. The inequality

$$0 = \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\mathbf{z}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \sum_{j=1}^q \mu_j h_j(\mathbf{z})$$

is inconsistent with at least one μ_j being positive and all $h_j(\mathbf{z})$ being negative. Hence, it suffices to assume all $\mu_j = 0$. For any $\mathbf{y} \in C$ we now find that

$$0 = \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \sum_{i=1}^p \lambda_i g_i(\mathbf{y}).$$

Let $a(\mathbf{y})$ denote the affine function $\sum_{i=1}^p \lambda_i g_i(\mathbf{y})$. We have just shown that $a(\mathbf{y}) \geq 0$ for all \mathbf{y} in a neighborhood of \mathbf{x} . Because $a(\mathbf{x}) = 0$, Fermat's rule requires the vector $\nabla a(\mathbf{x})$ to vanish. Finally, the fact that some of the λ_i are nonzero contradicts the assumed linear independence of the gradient vectors $\nabla g_i(\mathbf{x})$. The only possibility left is $\lambda_0 > 0$. Divide $\mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ by λ_0 to achieve the canonical form of the multiplier rule.

Finally for the converse, let $\mathbf{y} \in C$ be any feasible point. The inequalities

$$f(\mathbf{x}) = \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\mathbf{y})$$

follow directly from properties (a) through (c) and demonstrate that \mathbf{x} furnishes the constrained minimum of $f(\mathbf{y})$. ■

Proposition 14.7.1 falls short of our expectations in the sense that the usual multiplier rule involves a stationarity condition. For convex programs, the required gradients do not necessarily exist. Fortunately, subgradients provide a suitable substitute. One can better understand the situation by exploiting the fact that \mathbf{x} minimizes the Lagrangian over the convex set C . This suggests that we replace the Lagrangian by the related function

$$\mathcal{K}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \mathcal{L}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu}) + \delta_C(\mathbf{y}).$$

Because $\mathcal{K}(\mathbf{y}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ achieves its unconstrained minimum over \mathbb{R}^n at \mathbf{x} , we have $\mathbf{0} \in \partial \mathcal{K}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$. If the sum rule for subdifferentials is justified, then in view of Example 14.4.1 we recover the generalization

$$\mathbf{0} \in \partial f(\mathbf{x}) + \sum_{i=1}^p \lambda_i \partial g_i(\mathbf{x}) + \sum_{j=1}^q \mu_j \partial h_j(\mathbf{x}) + N_C(\mathbf{x})$$

of the usual multiplier rule. As we have previously stressed, simple sufficient conditions ensure the validity of the sum rule.

As an example, consider minimizing the function

$$f(\mathbf{x}) = \sum_i a_i |x_i| + \sum_i b_i x_i$$

over the closed unit ball, where each constant $a_i \geq 0$. The multiplier conditions are

$$0 \in b_i + \mu x_i + a_i \begin{cases} 1 & x_i > 0 \\ [-1, 1] & x_i = 0 \\ -1 & x_i < 0. \end{cases}$$

If all $|b_i| \leq a_i$, then these are satisfied at the origin $\mathbf{0}$ with the choice $\mu = 0$ dictated by complementary slackness. If any $|b_i| > a_i$, then we take μ to be the positive square root of

$$\mu^2 = \sum_{|b_i| > a_i} [b_i - a_i \operatorname{sgn}(b_i)]^2.$$

Those components x_i with $|b_i| \leq a_i$ we assign the value 0. The remaining components we assign the value

$$x_i = \frac{1}{\mu} [-b_i + a_i \operatorname{sgn}(b_i)].$$

Because the function $f(\mathbf{x})$ is homogeneous and its minimum is negative, the optimal point \mathbf{x} occurs on the boundary of the unit ball.

14.8 Problems

1. Derive the Fenchel conjugates displayed in Table 14.1 for functions on the real line.
2. Show that the function $f_1(x) = (x^2 - 1)^2$ has Fenchel biconjugate

$$f_1^{**}(x) = \begin{cases} 0 & |x| \leq 1 \\ (x^2 - 1)^2 & |x| > 1 \end{cases}$$

and that the function

$$f_2(x) = \begin{cases} |x| & |x| \leq 1 \\ 2 - |x| & 1 < |x| \leq 3/2 \\ |x| - 1 & |x| > 3/2 \end{cases}$$

has Fenchel biconjugate

$$f_2^{**}(x) = \begin{cases} \frac{|x|}{3} & |x| \leq 3/2 \\ f_2(x) & |x| > 3/2. \end{cases}$$

TABLE 14.1. Some specific Fenchel conjugates

$f(x)$	$f^*(y)$
x	$\begin{cases} 0 & y = 1 \\ \infty & y \neq 1 \end{cases}$
$ x $	$\begin{cases} 0 & y \leq 1 \\ \infty & y > 1 \end{cases}$
$\begin{cases} \frac{1}{x} & x > 0 \\ \infty & x \leq 0 \end{cases}$	$\begin{cases} \infty & y > 0 \\ -2\sqrt{-y} & y \leq 0 \end{cases}$
$\begin{cases} x \ln x & x > 0 \\ \infty & x \leq 0 \end{cases}$	e^{y-1}
e^x	$\begin{cases} y \ln y - y & y > 0 \\ 0 & y = 0 \\ \infty & y < 0 \end{cases}$
$\begin{cases} x \ln x + (1-x) \ln(1-x) & x \in (0, 1) \\ 0 & x \in \{0, 1\} \\ \infty & \text{otherwise} \end{cases}$	$\ln(1 + e^y)$

(Hints: There is no need to calculate the conjugate $f_i^*(y)$. According to inequality (14.2), $f_i^{**}(x)$ falls below $f_i(x)$. It also falls above any line supporting $f_i(x)$.)

3. Prove the Fenchel-Young inequality $xy \leq e^x - y + y \ln y$ for x and y nonnegative.
4. Find the support function of the convex set $C = \{\mathbf{x} : x_2 + e^{x_1} \leq 0\}$ in \mathbb{R}^2 . Explain why this is pertinent to Problem 23 of Chap. 6.
5. Show that the Fenchel conjugate of $f(\mathbf{x}) = \max_{1 \leq i \leq n} x_i$ is

$$f^*(\mathbf{y}) = \begin{cases} 0 & \text{all } y_i \geq 0 \text{ and } \sum_{i=1}^n y_i = 1 \\ \infty & \text{otherwise .} \end{cases}$$

6. Assume $f(x)$ is a continuous, strictly increasing function with functional inverse $g(y)$ and value $f(0) = 0$. Show that the even functions defined by

$$F(x) = \int_0^x f(u)du \quad \text{and} \quad G(y) = \int_0^y g(v)dv$$

for $x \geq 0$ and $y \geq 0$ constitute a conjugate pair. Why is Example 1.2.6 a special case? Give a graphical interpretation of Young's inequality $xy \leq F(x) + G(y)$. (Hint: Equality holds in Young's inequality when $x = g(y)$. Interpret graphically and in terms of $F^*(y)$.)

7. Suppose $f(\mathbf{x})$ is a convex function and $\mathbf{u} \in \partial f(\mathbf{x})$ and $\mathbf{v} \in \partial f(\mathbf{y})$. Prove that $(\mathbf{u} - \mathbf{v})^*(\mathbf{x} - \mathbf{y}) \geq 0$.
8. For some orthogonal matrix \mathbf{O} , suppose $f(\mathbf{x}) = f(\mathbf{O}\mathbf{x})$ for all \mathbf{x} . Prove that $f^*(\mathbf{y}) = f^*(\mathbf{O}\mathbf{y})$ for all \mathbf{y} as well.
9. Suppose the continuous function $f(\mathbf{x})$ defined on \mathbb{R}^n satisfies the condition $\lim_{\|\mathbf{x}\| \rightarrow \infty} \|\mathbf{x}\|^{-1} f(\mathbf{x}) = \infty$. Show that $f^*(\mathbf{y}) < \infty$ for all \mathbf{y} .
10. Demonstrate that $f(\mathbf{x}) = \frac{1}{2}\|\mathbf{x}\|^2$ is the only function satisfying the identity $f(\mathbf{x}) = f^*(\mathbf{x})$ for all \mathbf{x} . (Hint: Use the trivial inequality $f^*(\mathbf{y}) + f(\mathbf{x}) \geq \mathbf{y}^*\mathbf{x}$ to prove that $f(\mathbf{x}) \geq \frac{1}{2}\|\mathbf{x}\|^2$ when $f(\mathbf{x}) = f^*(\mathbf{x})$. For the reverse inequality, substitute this inequality in the definition of $f^*(\mathbf{y})$.)
11. Let $f(\mathbf{y})$ be a differentiable function from \mathbb{R}^n to \mathbb{R} . Prove that \mathbf{x} is a global minimum of $f(\mathbf{y})$ if and only if $\nabla f(\mathbf{x}) = \mathbf{0}$ and $f^{**}(\mathbf{x}) = f(\mathbf{x})$ [130]. (Hints: If \mathbf{x} is a global minimum, then $\mathbf{0}$ is a subgradient of $f(\mathbf{y})$ at \mathbf{x} . Conversely, if the two conditions hold, then show that every directional derivative $d_{\mathbf{v}}f(\mathbf{x})$ satisfies $d_{\mathbf{v}}f^{**}(\mathbf{x}) \leq 0$. Because $-d_{-\mathbf{v}}f^{**}(\mathbf{x}) \leq d_{\mathbf{v}}f^{**}(\mathbf{x})$, we have in fact $d_{\mathbf{v}}f^{**}(\mathbf{x}) = 0$ for every direction \mathbf{v} . Now use Problem 21 of Chap. 4 to establish that $\nabla f^{**}(\mathbf{x}) = \mathbf{0}$. Because $f^{**}(\mathbf{y})$ is convex, \mathbf{x} minimizes $f^{**}(\mathbf{y})$.)
12. Let the convex function $f(\mathbf{x}, \mathbf{y})$ have Fenchel conjugate $f^*(\mathbf{u}, \mathbf{v})$. Demonstrate that the function $g(\mathbf{x}) = \inf_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})$ has Fenchel conjugate $f^*(\mathbf{u}, \mathbf{0})$.
13. Let $B = \{\mathbf{x} : \|\mathbf{x}\|_{\dagger} \leq 1\}$ be the closed unit ball associated with a norm $\|\mathbf{x}\|_{\dagger}$ on \mathbb{R}^n . Show that $\|\mathbf{y}\|_{\star} = \delta_B^*(\mathbf{y}) = \sup_{\mathbf{x} \in B} \mathbf{y}^*\mathbf{x}$ also qualifies as a norm.
14. Assume $f(t)$ is a proper even function from \mathbb{R} to $(-\infty, \infty]$. Let $\|\mathbf{x}\|_{\dagger}$ be a norm on \mathbb{R}^n with dual norm $\|\mathbf{y}\|_{\star}$. Prove that the composite function $g(\mathbf{x}) = f(\|\mathbf{x}\|_{\dagger})$ has Fenchel conjugate $g^*(\mathbf{y}) = f^*(\|\mathbf{y}\|_{\star})$. We have already considered the special case of the self-conjugate function $f(t) = \frac{1}{2}t^2$. (Hint: $g^*(\mathbf{y}) = \sup_{t \geq 0} \sup_{\|\mathbf{x}\|_{\dagger} = t} [\mathbf{y}^*\mathbf{x} - f(t)]$.)
15. The infimal convolution of two convex functions $f(\mathbf{x})$ and $g(\mathbf{x})$ is defined by $(f \odot g)(\mathbf{x}) = \inf_{\mathbf{w}} [f(\mathbf{w}) + g(\mathbf{x} - \mathbf{w})]$. Prove that $(f \odot g)(\mathbf{x})$ is convex and has Fenchel conjugate $f^*(\mathbf{y}) + g^*(\mathbf{y})$. Calculate $(f \odot g)(\mathbf{x})$ when $f(\mathbf{x}) = \|\mathbf{x}\|$ and $g(\mathbf{x}) = \delta_U(\mathbf{y})$ for a nonempty set U . What is the Fenchel conjugate of this particular infimal convolution?
16. Suppose that the convex function $f(\mathbf{x})$ is coercive and twice continuously differentiable with $d^2f(\mathbf{x})$ positive definite for all \mathbf{x} . Argue via the implicit function theorem that the stationarity equation $\mathbf{0} = \mathbf{y} - \nabla f(\mathbf{x})$ can be solved for \mathbf{x} in terms of \mathbf{y} . Furthermore, show that

$\mathbf{x}(\mathbf{y}) = \nabla f^*(\mathbf{y})$ by differentiating the equation $f^*(\mathbf{y}) = \mathbf{y}^* \mathbf{x}(\mathbf{y}) - f[\mathbf{x}(\mathbf{y})]$. Finally, show that $d^2 f^*(\mathbf{y})$ exists and equals the matrix inverse of $d^2 f[\mathbf{x}(\mathbf{y})]$. (Hints: For the last assertion, consider the equation $\nabla f[\nabla f^*(\mathbf{y})] = \mathbf{y}$, and apply Problem 20 of Chap. 4. Section 12.3 defines coercive functions.)

17. Let $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x}$ for \mathbf{A} positive semidefinite but not positive definite. Prove that $f^*(\mathbf{A} \mathbf{z} + \mathbf{b}) = \frac{1}{2} \mathbf{z}^* \mathbf{A} \mathbf{z}$ and $f^*(\mathbf{y}) = \infty$ if \mathbf{y} cannot be represented as $\mathbf{A} \mathbf{z} + \mathbf{b}$ for some \mathbf{z} . (Hint: Invoke Problem 28 of Chap. 5.)
18. Let \mathbf{A} and \mathbf{B} be positive semidefinite matrices of the same dimension. Show that the matrix $a\mathbf{A} + b\mathbf{B}$ is positive semidefinite for every pair of nonnegative scalars a and b . Thus, the set of positive semidefinite matrices is a convex cone. Why is it a closed set as well?
19. Let \mathbf{A} and \mathbf{B} be positive definite matrices of the same dimension. We write $\mathbf{A} \succeq \mathbf{B}$ provided $\mathbf{x}^* \mathbf{A} \mathbf{x} \geq \mathbf{x}^* \mathbf{B} \mathbf{x}$ for all vectors \mathbf{x} . Demonstrate via the Fenchel conjugate that $\mathbf{A} \succeq \mathbf{B}$ implies $\mathbf{B}^{-1} \succeq \mathbf{A}^{-1}$.
20. Show that the two closed convex cones

$$\begin{aligned} C_1 &= \{ \mathbf{x} \in \mathbb{R}^n : x_i \geq 0, \forall i \} \\ C_2 &= \{ \mathbf{x} \in \mathbb{R}^n : x_i \leq x_{i+1}, \forall i < n \} \end{aligned}$$

have the polar cones

$$\begin{aligned} C_1^\circ &= \{ \mathbf{y} \in \mathbb{R}^n : y_i \leq 0, \forall i \} \\ C_2^\circ &= \{ \mathbf{y} \in \mathbb{R}^n : \sum_{i=1}^j y_i \geq 0, \forall j < n, \text{ and } \sum_{i=1}^n y_i = 0 \}. \end{aligned}$$

21. For a closed convex set C and $\mathbf{x} \in C$, prove that the normal cone

$$N_C(\mathbf{x}) = \{ \mathbf{y} : P_C(\mathbf{x} + \mathbf{y}) = \mathbf{x} \},$$

where $P_C(\mathbf{z})$ projects a point \mathbf{z} onto the closest point in C .

22. Let C be a closed convex cone in \mathbb{R}^n . Find the normal cone $N_C(\mathbf{x})$ when $\mathbf{x} \in C$.
23. Let C be a closed convex cone. For any point \mathbf{x} , verify the representation $\mathbf{x} = P_C(\mathbf{x}) + P_{C^\circ}(\mathbf{x})$ with $P_C(\mathbf{x})^* P_{C^\circ}(\mathbf{x}) = 0$. Here P_K denotes projection onto the closed convex set K . (Hints: Let $\mathbf{a} = P_C(\mathbf{x})$ and $\mathbf{b} = P_{C^\circ}(\mathbf{x})$. Prove that $(r\mathbf{a} - \mathbf{a})^*(\mathbf{x} - \mathbf{a}) \leq 0$ and $(r\mathbf{b} - \mathbf{b})^*(\mathbf{x} - \mathbf{b}) \leq 0$ for $r \in \{0, 2\}$. Next show that $\mathbf{x} - \mathbf{a} \in C^\circ$ and that $\mathbf{x} - \mathbf{b} \in C$. Finally, show that $\mathbf{x} - \mathbf{a} = \mathbf{b}$. Throughout apply the obtuse angle criterion of Example 6.5.3 and the definition of a polar cone.)

24. Suppose the $n \times n$ matrix \mathbf{A} is positive definite. Show that the function $\|\mathbf{x}\|_{\dagger} = \sqrt{\mathbf{x}^* \mathbf{A} \mathbf{x}}$ is a norm. Find its dual norm. (Hint: $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^* \mathbf{A} \mathbf{y}$ defines an inner product.)

25. A norm $\|\mathbf{x}\|_{\dagger}$ is said to be strictly convex if the conditions

$$\|\mathbf{x}\|_{\dagger} = \|\mathbf{y}\|_{\dagger} = \left\| \frac{1}{2}(\mathbf{x} + \mathbf{y}) \right\|_{\dagger} = 1$$

imply $\mathbf{x} = \mathbf{y}$. Which of the norms $\|\mathbf{x}\|_1$, $\|\mathbf{x}\|$, and $\|\mathbf{x}\|_{\infty}$ is strictly convex? When $\|\mathbf{x}\|_{\dagger}$ is strictly convex, prove that the closest point in a closed convex set C to an external point \mathbf{z} is unique. (Hint: Reduce the problem to the case $\mathbf{z} = \mathbf{0}$.)

26. Show that a closed convex function $f(\mathbf{x})$ satisfies

$$f(\mathbf{x}) = \sup_{r>0} \inf_{0<\|\mathbf{y}-\mathbf{x}\|\leq r} f(\mathbf{y})$$

provided $\text{dom } f$ contains at least two points. (Hint: Convexity gives a lower bound and lower semicontinuity an upper bound.)

27. Find the subdifferentials of the functions

$$f_1(x) = \begin{cases} 0 & x \in [-1, 1] \\ |x| - 1 & x \in [-2, -1) \cup (1, 2] \\ \infty & \text{otherwise} \end{cases}$$

$$f_2(x) = \begin{cases} 1 - \sqrt{1 - x^2} & x \in [-1, 1] \\ \infty & \text{otherwise.} \end{cases}$$

28. Use the representation $|x| = \max\{-x, x\}$ to find the subdifferential of $|x|$.

29. Calculate the subdifferentials $\partial\|\mathbf{x}\|_1$ and $\partial\|\mathbf{x}\|_{\infty}$ for $\mathbf{x} \in \mathbb{R}^2$ at the points $\mathbf{x} = (0, 0)$, $\mathbf{x} = (1, 0)$, and $\mathbf{x} = (1, 1)$.

30. Let $f(x)$ and $g(x)$ be positive increasing convex functions with common essential domain equal to an interval J . Prove that the product $f(x)g(x)$ is convex with subdifferential $f(x)\partial g(x) + g(x)\partial f(x)$ on the interior of J . (Hints: First prove convexity, and then take forward directional derivatives.)

31. Consider a positive concave function $f(x)$ with essential domain an interval J . Show that the function $g(x) = f(x)^{-1}$ is convex with subdifferential $f(x)^{-2}\partial[-f(x)]$ for x interior to J . (Hints: First prove convexity, and then take forward directional derivatives.)

32. Consider the indicator function $\delta_C(\mathbf{x})$ of the set

$$C = \{\mathbf{x} \in \mathbb{R}^2 : x_1^2 + (x_2 - 1)^2 \leq 1\}.$$

Prove that

$$\partial\delta_C(\mathbf{0}) = \{\mathbf{g} \in \mathbb{R}^2 : g_1 = 0, g_2 \leq 0\}$$

and that

$$d_{\mathbf{v}}\delta_C(\mathbf{0}) = \infty \neq 0 = \sup_{\mathbf{g} \in \partial\delta_C(\mathbf{0})} \mathbf{g}^* \mathbf{v}$$

for $\mathbf{v} = (1, 0)^*$. This result is inconsistent with the support function of $\partial\delta_C(\mathbf{0})$ attaining the upper bound $d_{\mathbf{v}}\delta_C(\mathbf{0})$. In this case $\mathbf{0}$ does not belong to the interior of $\text{dom}(\delta_C)$.

33. Demonstrate that the sum $C + D$ of a compact set C and a closed set D is closed.
34. Let $f(x)$ equal $-\sqrt{x}$ for $x \geq 0$ and ∞ for $x < 0$. If $g(x) = f(-x)$, then show that $\partial f(0) = \emptyset$ and $\partial g(0) = \emptyset$ but $\partial[f(0) + g(0)] = \mathbb{R}$. This result appears to contradict the sum rule. What assumption in our derivation of the sum rule fails?
35. A counterexample to the chain rule can be constructed by considering the closed convex set $C = \{\mathbf{x} \in \mathbb{R}^2 : x_1 x_2 \geq 1, x_1 > 0, x_2 > 0\}$. Show that the Fenchel conjugate of the indicator $\delta_C(\mathbf{y})$ equals the support function

$$\delta^*(\mathbf{y}) = \begin{cases} -2\sqrt{y_1 y_2} & y_1 \leq 0 \text{ and } y_2 \leq 0 \\ \infty & \text{otherwise.} \end{cases}$$

Given the symmetric matrix

$$P = \begin{pmatrix} 0 & 0 \\ 0 & 1 \end{pmatrix}$$

that projects a point \mathbf{y} onto the y_2 axis, prove that

$$\begin{aligned} \partial(\delta_C^* \circ P)(\mathbf{0}) &= \{\mathbf{x} : x_1 = 0, x_2 \geq 0\} \\ P^* \partial\delta_C^*(P\mathbf{0}) &= \{\mathbf{x} : x_1 = 0, x_2 > 0\}. \end{aligned}$$

Thus, the two sets differ by the presence and absence of $\mathbf{0}$.

36. Let $f(\mathbf{y})$ be a convex function and C be a closed convex set in \mathbb{R}^m . For $\mathbf{x} \in \text{dom}(f) \cap C$, prove the equivalence of the following three statements:
- (a) \mathbf{x} minimizes $f(\mathbf{y})$ on C .
 - (b) $d_{\mathbf{v}}f(\mathbf{x}) \geq 0$ for all directions $\mathbf{v} = \mathbf{y} - \mathbf{x}$ defined by $\mathbf{y} \in C$.
 - (c) $\mathbf{0} \in \partial[f(\mathbf{x}) + \delta_C(\mathbf{x})]$.

Furthermore, deduce that

$$\partial[f(\mathbf{x}) + \delta_C(\mathbf{x})] = \partial f(\mathbf{x}) + N_C(\mathbf{x})$$

when \mathbf{x} belongs to the interior of either $\text{dom}(f)$ or C .

37. Let $f(\mathbf{y}) : \mathbb{R}^n \mapsto \mathbb{R}$ be a convex differentiable function, \mathbf{A} an $n \times m$ matrix with j th column \mathbf{a}_j , and $\lambda > 0$ a constant. Prove the following assertions about the function $g(\mathbf{z}) = f(\mathbf{Az}) + \lambda \|\mathbf{z}\|_1$: (a) $g(\mathbf{z})$ achieves its global minimum at some point \mathbf{x} whenever $f(\mathbf{y})$ is bounded below or λ is sufficiently large, (b) at the point \mathbf{x}

$$\mathbf{a}_j^* \nabla f(\mathbf{Ax}) = \begin{cases} -\lambda & x_j > 0 \\ u \in [-\lambda, \lambda] & x_j = 0 \\ \lambda & x_j < 0 \end{cases}$$

for all j , and (c) for the choice $f(\mathbf{y}) = \frac{1}{2} \|\mathbf{u} - \mathbf{y}\|^2$ the components of the minimum point satisfy

$$x_j = \frac{S[\mathbf{a}_j^*(\mathbf{u} - \sum_{i \neq j} x_i \mathbf{a}_i), \lambda]}{\|\mathbf{a}_j\|^2},$$

where $S(v, \lambda) = \text{sgn}(v) \max\{|v| - \lambda, 0\}$ is the soft threshold operator.

38. Suppose the matrix \mathbf{X} has singular value decomposition $\mathbf{U}\mathbf{\Sigma}\mathbf{V}^*$ with all diagonal entries of $\mathbf{\Sigma}$ positive. Calculate the subdifferential

$$\partial \|\mathbf{X}\|_* = \{\mathbf{UV}^* + \mathbf{W} : \mathbf{U}^*\mathbf{W} = \mathbf{0}, \mathbf{W}\mathbf{V} = \mathbf{0}, \|\mathbf{W}\| \leq 1\}$$

of the nuclear norm [269]. (Hints: See Examples 14.3.6 and 14.4.2. The equality $\text{tr}(\mathbf{Y}^*\mathbf{X}) = \|\mathbf{X}\|_* \|\mathbf{Y}\|_2$ entails equality in Fan's inequality.)

39. Suppose the convex function $g(\mathbf{y} \mid \mathbf{x})$ majorizes the convex function $f(\mathbf{y})$ around the point \mathbf{x} . Demonstrate that $\partial f(\mathbf{x}) \subset \partial g(\mathbf{x} \mid \mathbf{x})$. Give an example where $\partial g(\mathbf{x} \mid \mathbf{x}) \neq \partial f(\mathbf{x})$. If $g(\mathbf{y} \mid \mathbf{x})$ is differentiable at \mathbf{x} , then equality holds.
40. The $\ell_{p,q}$ norm on \mathbb{R}^n is useful in group penalties [230]. Suppose the sets σ_g partition $\{1, \dots, n\}$. For $\mathbf{x} \in \mathbb{R}^n$ let \mathbf{x}_{σ_g} denote the vector formed by taking the components of \mathbf{x} derived from σ_g . For p and q between 1 and ∞ , the $\ell_{p,q}$ norm equals

$$\|\mathbf{x}\|_{p,q} = \left(\sum_g \|\mathbf{x}_{\sigma_g}\|_q^p \right)^{1/p},$$

Demonstrate that the $\ell_{r,s}$ norm is dual to the $\ell_{p,q}$ norm, where r and s satisfy $p^{-1} + r^{-1} = 1$ and $q^{-1} + s^{-1} = 1$.

15

Feasibility and Duality

15.1 Introduction

This chapter provides a concrete introduction to several advanced topics in optimization theory. Specifying an interior feasible point is the first issue that must be faced in applying a barrier method. Given an exterior point, Dykstra's algorithm [21, 70, 79] finds the closest point in the intersection $\cap_{i=0}^{r-1} C_i$ of a finite number of closed convex sets. If C_i is defined by the convex constraint $h_i(\mathbf{x}) \leq 0$, then one obvious tactic for finding an interior point is to replace C_i by the set $C_i(\epsilon) = \{\mathbf{x} : h_j(\mathbf{x}) \leq -\epsilon\}$ for some small $\epsilon > 0$. Projecting onto the intersection of the $C_i(\epsilon)$ then produces an interior point.

The method of alternating projections is faster in practice than Dykstra's algorithm, but it is only guaranteed to find some feasible point, not the closest feasible point to a given exterior point. Projection operators are specific examples of paracontractions. We study these briefly and their classical counterparts, contractions and strict contractions. Under the right hypotheses, a contraction T possesses a unique fixed point, and the sequence $\mathbf{x}_{m+1} = T(\mathbf{x}_m)$ converges to it regardless of the initial point \mathbf{x}_0 .

Duality is one of the deepest and most pervasive themes of modern optimization theory. It takes considerable mathematical maturity to appreciate this subtle topic, and it is impossible to do it justice in a short essay. Every convex program generates a corresponding dual program, which can be simpler to solve than the original or primal program. We show how to

construct dual programs and relate the absence of a duality gap to Slater's constraint qualification. We also point out important connections between duality and the Fenchel conjugate.

15.2 Dykstra's Algorithm

Example 2.5.2 demonstrates that the distance $\text{dist}(\mathbf{x}, C)$ from a point \mathbf{x} in \mathbb{R}^n to a set C is a uniformly continuous function of \mathbf{x} . If the set C is closed and convex, then $\text{dist}(\mathbf{x}, C)$ is convex in \mathbf{x} and $\text{dist}(\mathbf{x}, C) = \|P_C(\mathbf{x}) - \mathbf{x}\|$ for exactly one projected point $P_C(\mathbf{x}) \in C$. (See Examples 2.5.5 and 6.3.3 and Proposition 6.2.2.) Although it is generally impossible to calculate $P_C(\mathbf{x})$ explicitly, some specific projection operators are well known.

Example 15.2.1 *Examples of Projection Operators*

Closed Euclidean Ball: If $C = \{\mathbf{y} \in \mathbb{R}^n : \|\mathbf{y} - \mathbf{z}\| \leq r\}$, then

$$P_C(\mathbf{x}) = \begin{cases} \mathbf{z} + \frac{r}{\|\mathbf{x} - \mathbf{z}\|}(\mathbf{x} - \mathbf{z}) & \mathbf{x} \notin C \\ \mathbf{x} & \mathbf{x} \in C \end{cases}.$$

Closed Rectangle: If $C = [a, b]$ is a closed rectangle in \mathbb{R}^n , then

$$P_C(\mathbf{x})_i = \begin{cases} a_i & x_i < a_i \\ x_i & x_i \in [a_i, b_i] \\ b_i & x_i > b_i \end{cases}.$$

Hyperplane: If $C = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{a}^* \mathbf{y} = b\}$ for $\mathbf{a} \neq \mathbf{0}$, then

$$P_C(\mathbf{x}) = \mathbf{x} - \frac{\mathbf{a}^* \mathbf{x} - b}{\|\mathbf{a}\|^2} \mathbf{a}.$$

Closed Halfspace: If $C = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{a}^* \mathbf{y} \leq b\}$ for $\mathbf{a} \neq \mathbf{0}$, then

$$P_C(\mathbf{x}) = \begin{cases} \mathbf{x} - \frac{\mathbf{a}^* \mathbf{x} - b}{\|\mathbf{a}\|^2} \mathbf{a} & \mathbf{a}^* \mathbf{x} > b \\ \mathbf{x} & \mathbf{a}^* \mathbf{x} \leq b \end{cases}.$$

Subspace: If C is the range of a matrix \mathbf{A} with full column rank, then

$$P_C(\mathbf{x}) = \mathbf{A}(\mathbf{A}^* \mathbf{A})^{-1} \mathbf{A}^* \mathbf{x}.$$

Positive Semidefinite Matrices: Let \mathbf{M} be an $n \times n$ symmetric matrix with spectral decomposition $\mathbf{M} = \mathbf{U} \mathbf{D} \mathbf{U}^*$, where \mathbf{U} is an orthogonal matrix and \mathbf{D} is a diagonal matrix with i th diagonal entry d_i . The projection of \mathbf{M} onto the set S of positive semidefinite matrices is given by $P_S(\mathbf{M}) = \mathbf{U} \mathbf{D}_+ \mathbf{U}^*$, where \mathbf{D}_+ is diagonal with i th diagonal entry $\max\{d_i, 0\}$. Problem 7 asks the reader to check this fact.

Unit Simplex: The problem of projecting a point \mathbf{x} onto the unit simplex

$$S = \left\{ \mathbf{y} : \sum_{i=1}^n y_i = 1, y_i \geq 0 \text{ for } 1 \leq i \leq n \right\}.$$

can be solved by a simple algorithm [77]. If we apply Gibbs lemma as sketched in Problem 18 of Chap. 5 to $f(\mathbf{y}) = \frac{1}{2} \|\mathbf{y} - \mathbf{x}\|^2$, then the coordinates of the minimum point \mathbf{y} satisfy

$$y_i = \begin{cases} x_i - \lambda & y_i > 0 \\ x_i - \lambda + \mu_i & y_i = 0 \end{cases}$$

for Lagrange multipliers λ and $\mu_i \geq 0$. Setting $I_+ = \{i : y_i > 0\}$ and invoking the equality constraint

$$1 = \sum_{i \in I_+} y_i = \sum_{i \in I_+} x_i - |I_+| \lambda$$

then imply

$$\lambda = \frac{1}{|I_+|} \left(\sum_{i \in I_+} x_i - 1 \right).$$

The catch, of course, is that we do not know I_+ . The key to avoiding all 2^n possible subsets is the simple observation that the y_i and x_i are consistently ordered. Suppose on the contrary that $x_i < x_j$ and $y_j < y_i$. For small $s > 0$ substitute $y_j + s$ for y_j and $y_i - s$ for y_i . The objective function $f(\mathbf{y})$ then changes by the amount

$$\begin{aligned} & \frac{1}{2} \left[(y_i - s - x_i)^2 + (y_j + s - x_j)^2 - (y_i - x_i)^2 - (y_j - x_j)^2 \right] \\ &= s(x_i - x_j + y_j - y_i) + s^2, \end{aligned}$$

which is negative for small s . Thus, let z_1, \dots, z_n denote the x_i ordered from largest to smallest. For each integer j between 1 and n , it suffices to set $\lambda = \frac{1}{j} (\sum_{i=1}^j z_i - 1)$ and check whether $z_j > \lambda$ and $z_{j+1} \leq \lambda$. When these two conditions are met, we put $y_i = (x_i - \lambda)_+$ for all i . Michelot's [196] algorithm as described in Problem 13 also solves this projection problem.

Closed ℓ_1 Ball: Projecting a point \mathbf{x} onto the ℓ_1 ball

$$C = \{ \mathbf{y} : \|\mathbf{y} - \mathbf{z}\|_1 \leq r \}.$$

yields to a variation of the previous projection algorithm as sketched in Problem 14.

Isotone Convex Cone: Projection onto the convex cone

$$C = \{\mathbf{y} : y_i \leq y_{i+1} \text{ for } 1 \leq i < n\}.$$

is accomplished by the pool adjacent violators algorithm as discussed later in Example 15.2.3. ■

Dykstra's algorithm [21, 70, 79] is designed to find the projection of a point \mathbf{x} onto a finite intersection $C = \bigcap_{i=0}^{r-1} C_i$ of r closed convex sets. Here are some possible situations where Dykstra's algorithm applies.

Example 15.2.2 *Applications of Dykstra's Algorithm*

Linear Equalities: Any solution of the system of linear equations $\mathbf{Ax} = \mathbf{b}$ belongs to the intersection of the hyperplanes $\mathbf{a}_i^* \mathbf{x} = b_i$, where \mathbf{a}_i^* is the i th row of \mathbf{A} .

Linear Inequalities: Any solution of the system of linear inequalities $\mathbf{Ax} \leq \mathbf{b}$ belongs to the intersection of the halfspaces $\mathbf{a}_i^* \mathbf{x} \leq b_i$, where \mathbf{a}_i^* is the i th row of \mathbf{A} .

Isotone Regression: The least squares problem of minimizing the sum $\sum_{i=1}^n (x_i - w_i)^2$ subject to the constraints $w_i \leq w_{i+1}$ corresponds to projection of \mathbf{x} onto the intersection of the halfspaces

$$C_i = \{\mathbf{w} \in \mathbb{R}^n : w_i - w_{i+1} \leq 0\}, \quad 1 \leq i \leq n-1.$$

Convex Regression: The least squares problem of minimizing the sum $\sum_{i=1}^n (x_i - w_i)^2$ subject to the constraints $w_i \leq \frac{1}{2}(w_{i-1} + w_{i+1})$ corresponds to projection of \mathbf{x} onto the intersection of the halfspaces

$$C_i = \left\{ \mathbf{w} \in \mathbb{R}^n : w_i - \frac{1}{2}(w_{i-1} + w_{i+1}) \leq 0 \right\}, \quad 2 \leq i \leq n-1.$$

Quadratic Programming: To minimize the strictly convex quadratic form $\frac{1}{2} \mathbf{x}^* \mathbf{Ax} + \mathbf{b}^* \mathbf{x} + c$ subject to $\mathbf{Dx} = \mathbf{e}$ and $\mathbf{Fx} \leq \mathbf{g}$, we make the change of variables $\mathbf{y} = \mathbf{A}^{1/2} \mathbf{x}$. This transforms the problem to one of minimizing

$$\begin{aligned} \frac{1}{2} \mathbf{x}^* \mathbf{Ax} + \mathbf{b}^* \mathbf{x} + c &= \frac{1}{2} \|\mathbf{y}\|^2 + \mathbf{b}^* \mathbf{A}^{-1/2} \mathbf{y} + c \\ &= \frac{1}{2} \|\mathbf{y} + \mathbf{A}^{-1/2} \mathbf{b}\|^2 - \frac{1}{2} \mathbf{b}^* \mathbf{A}^{-1} \mathbf{b} + c \end{aligned}$$

subject to $\mathbf{DA}^{-1/2} \mathbf{y} = \mathbf{e}$ and $\mathbf{FA}^{-1/2} \mathbf{y} \leq \mathbf{g}$. The solution in the \mathbf{y} coordinates is determined by projecting $-\mathbf{A}^{-1/2} \mathbf{b}$ onto the convex feasible region determined by the revised constraints. Instead of the symmetric square root transformation $\mathbf{y} = \mathbf{A}^{1/2} \mathbf{x}$, one can employ the asymmetric square root transformation $\mathbf{y} = \mathbf{Ux}$ furnished by the Cholesky decomposition $\mathbf{L} = \mathbf{U}^*$ of \mathbf{A} . ■

To state Dykstra's algorithm, it is helpful to label the closed convex sets C_0, \dots, C_{r-1} and denote their intersection by $C = \bigcap_{i=0}^{r-1} C_i$. The algorithm keeps track of a primary sequence \mathbf{x}_m and a companion sequence $\boldsymbol{\omega}_m$. In the limit, \mathbf{x}_m tends to $P_C(\mathbf{x})$. To initiate the process, we set $\mathbf{x}_0 = \mathbf{x}$ and $\boldsymbol{\omega}_{-r+1} = \dots = \boldsymbol{\omega}_0 = \mathbf{0}$. For $m \geq 0$ we then iterate via

$$\begin{aligned}\mathbf{x}_{m+1} &= P_{C_{m \bmod r}}(\mathbf{x}_m + \boldsymbol{\omega}_{m-r+1}) \\ \boldsymbol{\omega}_{m+1} &= \mathbf{x}_m + \boldsymbol{\omega}_{m-r+1} - \mathbf{x}_{m+1}.\end{aligned}$$

Here $m \bmod r$ is the nonnegative remainder after dividing m by r . In essence, the algorithm cycles through the convex sets and projects the sum of the current vector and the relevant previous companion vector onto the current convex set.

TABLE 15.1. Iterates of Dykstra's algorithm

Iteration m	x_{m1}	x_{m2}
0	-1.00000	2.00000
1	-0.44721	0.89443
2	0.00000	0.89443
3	-0.26640	0.96386
4	0.00000	0.96386
5	-0.14175	0.98990
10	0.00000	0.99934
15	-0.00454	0.99999
25	-0.00014	1.00000
30	0.00000	1.00000
35	0.00000	1.00000

As an example, suppose $r = 2$, C_0 is the closed unit ball in \mathbb{R}^2 , and C_1 is the closed halfspace with $x_1 \geq 0$. The intersection C is the right half ball centered at the origin. Table 15.1 records the iterates of Dykstra's algorithm starting from the point $\mathbf{x}_0 = (-1, 2)^*$ and their eventual convergence to the geometrically obvious solution $(0, 1)^*$.

When C_i is a subspace, Dykstra's algorithm can dispense with the corresponding companion subsequence $\boldsymbol{\omega}_m$. In this case, $\boldsymbol{\omega}_{m+1}$ is perpendicular to C_i whenever $m \bmod r = i$. Indeed, since $P_{C_i}(\mathbf{y})$ is a linear projection, we have

$$\begin{aligned}\mathbf{x}_{m+1} &= P_{C_i}(\mathbf{x}_m + \boldsymbol{\omega}_{m-r+1}) \\ &= P_{C_i}(\mathbf{x}_m) + P_{C_i}(\boldsymbol{\omega}_{m-r+1}) \\ &= P_{C_i}(\mathbf{x}_m)\end{aligned}$$

under the perpendicularity assumption. The initial condition $\omega_{i-r} = \mathbf{0}$, the identity

$$\begin{aligned}\omega_{m+1} &= \mathbf{x}_m - \mathbf{x}_{m+1} + \omega_{m-r+1} \\ &= \mathbf{x}_m - P_{C_i}(\mathbf{x}_m) + \omega_{m-r+1} \\ &= P_{C_i^\perp}(\mathbf{x}_m) + \omega_{m-r+1},\end{aligned}$$

and induction show that ω_{m+1} belongs to the perpendicular complement C_i^\perp if $m \bmod r = i$. When all of the C_i are subspaces, Dykstra's algorithm reduces to the method of alternating projections first studied by von Neumann. Example 15.5.6 will justify Dykstra's algorithm theoretically.

Example 15.2.3 *Pool Adjacent Violators Algorithm*

Dykstra's algorithm can be beat in specific problems. Consider weighted isotone regression with objective function $f(\mathbf{x}) = \frac{1}{2} \sum_{i=1}^n w_i (y_i - x_i)^2$ and arbitrary positive weights w_i . The Lagrangian for this problem reads

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\mu}) = \frac{1}{2} \sum_{i=1}^n w_i (y_i - x_i)^2 + \sum_{i=1}^{n-1} \mu_i (x_i - x_{i+1}),$$

for nonnegative multipliers μ_i . The optimal point is determined by the stationarity conditions

$$0 = w_i (x_i - y_i) + \mu_i - \mu_{i-1}$$

and the complementary slackness conditions $\mu_i (x_i - x_{i+1}) = 0$. Here we define $\mu_0 = \mu_n = 0$ for convenience. Because of telescoping, one can solve for the multipliers in either of the equivalent forms

$$\mu_j = \sum_{i=1}^j w_i (y_i - x_i) = - \sum_{i=j+1}^n w_i (y_i - x_i). \quad (15.1)$$

since $\mu_n = \sum_{i=1}^n w_i (y_i - x_i) = 0$.

The pool adjacent violators algorithm [6, 157] exploits the fact that whole blocks of the solution vector \mathbf{x} are constant. From one block to the next, this constant increases. The algorithm starts with $x_i = y_i$ and $\mu_i = 0$ for all i and all blocks reducing to singletons. It then marches through the blocks and considers whether to consolidate or pool adjacent blocks. Let

$$B = \{l, l+1, \dots, r-1, r\}$$

denote a generic block. In view of complementary slackness, the right multiplier μ_r is assumed to be 0. According to the above calculations, the constant assigned to block B is the weighted average

$$x_B = \frac{1}{\sum_{i=l}^r w_i} \sum_{i=l}^r w_i y_i.$$

Now the only thing that prevents a block decomposition from supplying the solution vector is a reversal of two adjacent block constants. Suppose $B_1 = \{l_1, \dots, r_1\}$ and $B_2 = \{l_2, \dots, r_2\}$ are adjacent violating blocks. We pool the two blocks and assign the constant

$$x_{B_1 \cup B_2} = \frac{1}{\sum_{i=l_1}^{r_2} w_i} \sum_{i=l_1}^{r_2} w_i y_i$$

to $B_1 \cup B_2$. Note here that $r_1 + 1 = l_2$. We must also recalculate the multipliers associated with the combined block and check that they are nonnegative. Equation (15.1) is instrumental in achieving this goal. By induction we assume that it holds when 1 is replaced by l_i and n replaced by r_i for $i \in \{1, 2\}$. In recalculating the μ_j for j between l_1 and r_2 , we can take the last multiplier μ_{r_2} to be 0 by virtue of the definition of $x_{B_1 \cup B_2}$. For j between l_1 and r_1 , the left definition in equation (15.1) and the inequality $x_{B_1} > x_{B_1 \cup B_2}$ imply that the new $\mu_j \geq 0$. For j between l_2 and $r_2 - 1$, the right definition in equation (15.1) and the inequality $x_{B_1 \cup B_2} > x_{B_2}$ again imply that the new $\mu_j \geq 0$. Thus, the pooled block satisfies the multiplier conditions. Pooling continues until all blocks coalesce or no violations occur. At that point the multiplier conditions hold, and the minimum has been reached. ■

15.3 Contractive Maps

We met locally contractive maps in our study of convergence in Chap. 12. Here we discuss a generalization that is helpful in finding initial feasible points in nonlinear programming. First recall that a map $T : D \subset \mathbb{R}^n \mapsto \mathbb{R}^n$ is contractive relative to a norm $\|\mathbf{x}\|_{\dagger}$ if $\|T(\mathbf{y}) - T(\mathbf{z})\|_{\dagger} < \|\mathbf{y} - \mathbf{z}\|_{\dagger}$ for all $\mathbf{y} \neq \mathbf{z}$ in D . It is strictly contractive if there exists a constant $c \in [0, 1)$ with $\|T(\mathbf{y}) - T(\mathbf{z})\|_{\dagger} \leq c\|\mathbf{y} - \mathbf{z}\|_{\dagger}$ for all such pairs. Finally, T is said to be paracontractive provided for every fixed point \mathbf{y} of $T(\mathbf{x})$ the inequality $\|T(\mathbf{x}) - \mathbf{y}\|_{\dagger} < \|\mathbf{x} - \mathbf{y}\|_{\dagger}$ holds unless \mathbf{x} is itself a fixed point. A strictly contractive map is contractive, and a contractive map is paracontractive.

For instance, the affine map $\mathbf{M}\mathbf{x} + \mathbf{v}$ is contractive under some induced matrix norm whenever the spectral radius $\rho(\mathbf{M}) < 1$. (See Proposition 6.3.2 of the reference [166].) Projection onto a closed convex set C containing more than a single point is paracontractive but not contractive under the standard Euclidean norm. To validate paracontraction, note that the obtuse angle criterion stated in Example 6.5.3 implies

$$\begin{aligned} [\mathbf{x} - P_C(\mathbf{x})]^* [P_C(\mathbf{y}) - P_C(\mathbf{x})] &\leq 0 \\ [\mathbf{y} - P_C(\mathbf{y})]^* [P_C(\mathbf{x}) - P_C(\mathbf{y})] &\leq 0. \end{aligned}$$

Adding these inequalities, rearranging, and applying the Cauchy-Schwarz inequality give

$$\begin{aligned} \|P_C(\mathbf{x}) - P_C(\mathbf{y})\|^2 &\leq (\mathbf{x} - \mathbf{y})^* [P_C(\mathbf{x}) - P_C(\mathbf{y})] & (15.2) \\ &\leq \|\mathbf{x} - \mathbf{y}\| \cdot \|P_C(\mathbf{x}) - P_C(\mathbf{y})\|. \end{aligned}$$

Dividing the extremes of inequality (15.2) by $\|P_C(\mathbf{x}) - P_C(\mathbf{y})\|$ now demonstrates that $\|P_C(\mathbf{x}) - P_C(\mathbf{y})\| \leq \|\mathbf{x} - \mathbf{y}\|$. Equality holds in the Cauchy-Schwarz half of inequality (15.2) if and only if $P_C(\mathbf{x}) - P_C(\mathbf{y}) = c(\mathbf{x} - \mathbf{y})$ for some constant c . If \mathbf{y} is a fixed point, then overall equality in inequality (15.2) entails

$$\|P_C(\mathbf{x}) - P_C(\mathbf{y})\|^2 = (\mathbf{x} - \mathbf{y})^* [P_C(\mathbf{x}) - P_C(\mathbf{y})] = c\|\mathbf{x} - \mathbf{y}\|^2.$$

This precipitates a cascade of deductions. First that $c = 1$, second that $P_C(\mathbf{x}) - \mathbf{x} = P_C(\mathbf{y}) - \mathbf{y} = \mathbf{0}$, third that \mathbf{x} is a fixed point, and fourth that the projection map $P_C(\mathbf{x})$ is paracontractive under the Euclidean norm.

With this background in mind, we now state and prove an important result due to Elsner, Koltracht, and Neumann [85].

Proposition 15.3.1 *Suppose the continuous maps T_0, \dots, T_{r-1} of a set into itself are paracontractive under the norm $\|\mathbf{x}\|_{\dagger}$. Let F_i denote the set of fixed points of T_i . If the intersection $F = \bigcap_{i=0}^{r-1} F_i$ is nonempty, then the sequence*

$$\mathbf{x}_{m+1} = T_{m \bmod r}(\mathbf{x}_m)$$

converges to a limit in F . In particular, if $r = 1$ and $T = T_0$ has a nonempty set of fixed points F , then $\mathbf{x}_{m+1} = T(\mathbf{x}_m)$ converges to a point in F .

Proof: Let \mathbf{y} be any point in F . The scalar sequence $\|\mathbf{x}_m - \mathbf{y}\|_{\dagger}$ satisfies

$$\|\mathbf{x}_{m+1} - \mathbf{y}\|_{\dagger} = \|T_{m \bmod r}(\mathbf{x}_m) - \mathbf{y}\|_{\dagger} \leq \|\mathbf{x}_m - \mathbf{y}\|_{\dagger}$$

and therefore possesses a limit $d \geq 0$. Because the sequence \mathbf{x}_m is bounded, it possesses a cluster point \mathbf{x}_{∞} . Furthermore, $\|\mathbf{x}_{\infty} - \mathbf{y}\|_{\dagger}$ attains the lower bound d . But this implies that $\|T_i(\mathbf{x}_{\infty}) - \mathbf{y}\|_{\dagger} = \|\mathbf{x}_{\infty} - \mathbf{y}\|_{\dagger}$ for all i , and therefore $\mathbf{x}_{\infty} \in F$. For the choice $\mathbf{y} = \mathbf{x}_{\infty}$, the corresponding constant d equals 0. Finally, the monotone convergence of $\|\mathbf{x}_m - \mathbf{x}_{\infty}\|_{\dagger}$ to 0 implies $\lim_{m \rightarrow \infty} \mathbf{x}_m = \mathbf{x}_{\infty}$. ■

There are two corollaries to Proposition 15.3.1. First, the set of fixed points of the composite map $S = T_{r-1} \circ \dots \circ T_0$ equals $F = \bigcap_{i=0}^{r-1} F_i$. Second, S itself is paracontractive. Problem 21 asks the reader to prove these facts. As a trivial application of the proposition, consider the toy example of projection onto the half ball appearing in Table 15.1. For the chosen initial point \mathbf{x}_0 , a single round $P_{C_1} \circ P_{C_0}(\mathbf{x}_0)$ of projection lands in the half ball. In contrast to Dykstra's algorithm, the limit is not the closest point in $C = C_0 \cap C_1$ to \mathbf{x}_0 .

Alternatively, one can iterate via a convex combination

$$R(\mathbf{x}) = \sum_{i=0}^{r-1} \lambda_i T_i(\mathbf{x})$$

with positive coefficients. Clearly, any point $\mathbf{x} \in F$ is also a fixed point of $R(\mathbf{x})$. Conversely, if \mathbf{x} is a fixed point and \mathbf{y} is any point in F , then

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|_{\dagger} &= \left\| \sum_{i=0}^{r-1} \lambda_i T_i(\mathbf{x}) - \sum_{i=0}^{r-1} \lambda_i T_i(\mathbf{y}) \right\|_{\dagger} \\ &\leq \sum_{i=0}^{r-1} \lambda_i \|T_i(\mathbf{x}) - T_i(\mathbf{y})\|_{\dagger} \\ &\leq \sum_{i=0}^{r-1} \lambda_i \|\mathbf{x} - \mathbf{y}\|_{\dagger} \\ &= \|\mathbf{x} - \mathbf{y}\|_{\dagger}. \end{aligned}$$

Because equality must occur throughout, it follows from the paracontractiveness of T_i that $\mathbf{x} \in F_i$ for each i . Hence, the fixed points of $R(\mathbf{x})$ coincide with F . If we take $\mathbf{x} \notin F$ and $\mathbf{y} \in F$, then basically the same argument demonstrates the paracontractiveness requirement $\|R(\mathbf{x}) - \mathbf{y}\|_{\dagger} < \|\mathbf{x} - \mathbf{y}\|_{\dagger}$. Convergence of the iterates $\mathbf{x}_{n+1} = \sum_{i=0}^{r-1} \lambda_i T_i(\mathbf{x}_n)$ is now a consequence of the paracontractiveness of the map $R(\mathbf{x})$. The evidence suggests that simultaneous projection converges more slowly than alternating projection [42, 109]. However, simultaneous projection enjoys the advantage of being parallelizable.

Proposition 15.3.1 postulates the existence of fixed points. The next proposition introduces simple sufficient conditions guaranteeing existence.

Proposition 15.3.2 *Suppose T is contractive under a norm $\|\mathbf{x}\|_{\dagger}$ and maps the nonempty compact set $D \subset \mathbb{R}^n$ into itself. Then T has a unique fixed point in D . If T is a strict contraction with contraction constant c , then the assumption that D is compact can be relaxed to the assumption that D is closed.*

Proof: We first demonstrate that there is at most one fixed point \mathbf{y} . If there is a second fixed point $\mathbf{z} \neq \mathbf{y}$, then

$$\|\mathbf{y} - \mathbf{z}\|_{\dagger} = \|T(\mathbf{y}) - T(\mathbf{z})\|_{\dagger} < \|\mathbf{y} - \mathbf{z}\|_{\dagger},$$

which is a contradiction.

Now define $d = \inf_{\mathbf{x} \in D} f(\mathbf{x})$ for the function $f(\mathbf{x}) = \|T(\mathbf{x}) - \mathbf{x}\|_{\dagger}$. Since $f(\mathbf{x})$ is continuous, its infimum is attained at some point \mathbf{y} in the compact set D . If \mathbf{y} is not a fixed point, then

$$\|T \circ T(\mathbf{y}) - T(\mathbf{y})\|_{\dagger} < \|T(\mathbf{y}) - \mathbf{y}\|_{\dagger},$$

contradicting the definition of \mathbf{y} . On the other hand, if D is not compact, but T is a strict contraction, then choose any point $\mathbf{y} \in D$, and define the set $C = \{\mathbf{x} \in D : f(\mathbf{x}) \leq f(\mathbf{y})\}$. For $\mathbf{x} \in C$ we then have

$$\begin{aligned} \|\mathbf{x} - \mathbf{y}\|_{\dagger} &\leq \|\mathbf{x} - T(\mathbf{x})\|_{\dagger} + \|T(\mathbf{x}) - T(\mathbf{y})\|_{\dagger} + \|T(\mathbf{y}) - \mathbf{y}\|_{\dagger} \\ &\leq 2f(\mathbf{y}) + c\|\mathbf{x} - \mathbf{y}\|_{\dagger}. \end{aligned}$$

It follows that

$$\|\mathbf{x} - \mathbf{y}\|_{\dagger} \leq \frac{2f(\mathbf{y})}{1-c}.$$

Thus, the closed set C is bounded and hence compact. Furthermore, the inequality $f[T(\mathbf{x})] \leq cf(\mathbf{x})$ indicates that T maps C into itself. The rest of the argument proceeds as before. ■

Example 15.3.1 Stationary Distribution of a Markov Chain

Example 6.2.1 demonstrates that every finite state Markov chain possesses a stationary distribution. Under an appropriate ergodic hypothesis, this distribution is unique, and the chain converges to it. In understanding these phenomena, it simplifies notation to pass to column vectors and replace $\mathbf{P} = (p_{ij})$ by its transpose $\mathbf{Q} = (q_{ij})$. It is easy to check that \mathbf{Q} maps the standard simplex

$$S = \left\{ \mathbf{x} : x_i \geq 0, i = 1, \dots, n, \sum_{i=1}^n x_i = 1 \right\}$$

into itself and that candidate vectors belong to S . The natural norm on S is $\|\mathbf{x}\|_1 = \sum_{i=1}^n |x_i|$. According to Problem 9 of Chap. 2, the corresponding induced matrix norm of \mathbf{Q} is

$$\|\mathbf{Q}\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |q_{ij}| = \max_{1 \leq j \leq n} \sum_{i=1}^n p_{ji} = 1.$$

Thus, \mathbf{Q} is almost a contraction.

The standard ergodic hypothesis says that some power \mathbf{P}^k of \mathbf{P} has all entries positive. If we define $\mathbf{R} = \mathbf{Q}^k - c\mathbf{1}\mathbf{1}^*$ for c the minimum entry of \mathbf{P}^k , then the matrix \mathbf{R} has all entries nonnegative and norm $\|\mathbf{R}\|_1 < 1$. Since $\mathbf{1}^*(\mathbf{y} - \mathbf{x}) = 0$ for all pairs \mathbf{x} and \mathbf{y} in S , we have

$$\|\mathbf{Q}^k \mathbf{x} - \mathbf{Q}^k \mathbf{y}\|_1 = \|\mathbf{R}(\mathbf{x} - \mathbf{y})\|_1 \leq \|\mathbf{R}\|_1 \|\mathbf{x} - \mathbf{y}\|_1.$$

In other words, the map $\mathbf{x} \mapsto \mathbf{Q}^k \mathbf{x}$ is strictly contractive on S . It therefore possesses a unique fixed point \mathbf{y} , and $\lim_{m \rightarrow \infty} \mathbf{Q}^{mk} \mathbf{x} = \mathbf{y}$. Problem 22 asks the reader to check that the map $\mathbf{x} \mapsto \mathbf{Q} \mathbf{x}$ shares the fixed point \mathbf{y} and that $\lim_{m \rightarrow \infty} \mathbf{Q}^m \mathbf{x} = \mathbf{y}$. ■

Proposition 15.3.2 finds applications in many branches of mathematics and statistics. Part of its value stems from the guaranteed geometric rate of convergence of the iterates $\mathbf{x}_{m+1} = T(\mathbf{x}_m)$ to the fixed point \mathbf{y} . This assertion is a straightforward consequence of the inequality

$$\|\mathbf{x} - \mathbf{y}\|_{\dagger} \leq \|\mathbf{x} - T(\mathbf{x})\|_{\dagger} + \|T(\mathbf{x}) - T(\mathbf{y})\|_{\dagger} \leq f(\mathbf{x}) + c\|\mathbf{x} - \mathbf{y}\|_{\dagger}$$

involving the function $f(\mathbf{x}) = \|T(\mathbf{x}) - \mathbf{x}\|_{\dagger}$. It follows that

$$\|\mathbf{x} - \mathbf{y}\|_{\dagger} \leq \frac{1}{1-c}f(\mathbf{x}). \quad (15.3)$$

Substituting \mathbf{x}_m for \mathbf{x} and applying the inequality $f[T(\mathbf{x})] \leq cf(\mathbf{x})$ repeatedly now yield the geometric bound

$$\|\mathbf{x}_m - \mathbf{y}\|_{\dagger} \leq \frac{c^m}{1-c}f(\mathbf{x}_0).$$

In numerical practice, inequality (15.3) gives the preferred test for declaring convergence. If one wants \mathbf{x}_m to be within ϵ of the fixed point, then iteration should continue until $f(\mathbf{x}_m) \leq \epsilon(1-c)$.

15.4 Dual Functions

The Lagrange multiplier rule summarizes much of what we know about minimizing $f(\mathbf{x})$ subject to the constraints $g_i(\mathbf{x}) = 0$ for $1 \leq i \leq p$ and $h_j(\mathbf{x}) \leq 0$ for $1 \leq j \leq q$. Consequently, it is worth considering the standard Lagrangian function

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = f(\mathbf{x}) + \sum_{i=1}^p \lambda_i g_i(\mathbf{x}) + \sum_{j=1}^q \mu_j h_j(\mathbf{x})$$

in more detail. Here the multiplier vectors $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ are taken as arguments in addition to the variable \mathbf{x} . For a convex program satisfying a constraint qualification such as Slater's condition, a constrained global minimum $\hat{\mathbf{x}}$ of $f(\mathbf{x})$ is also an unconstrained global minimum of $\mathcal{L}(\mathbf{x}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}})$, where $\hat{\boldsymbol{\lambda}}$ and $\hat{\boldsymbol{\mu}}$ are the corresponding Lagrange multipliers. This fact is the content of Proposition 14.7.1

The behavior of $\mathcal{L}(\hat{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu})$ as a function of $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ is also interesting. Because $g_i(\hat{\mathbf{x}}) = 0$ for all i and $\hat{\boldsymbol{\mu}}_j h_j(\hat{\mathbf{x}}) = 0$ for all j , we have

$$\begin{aligned} \mathcal{L}(\hat{\mathbf{x}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}) - \mathcal{L}(\hat{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= \sum_{j=1}^q (\hat{\mu}_j - \mu_j) h_j(\hat{\mathbf{x}}) \\ &= - \sum_{j=1}^q \mu_j h_j(\hat{\mathbf{x}}) \\ &\geq 0. \end{aligned}$$

This proves the left inequality of the two saddle point inequalities

$$\mathcal{L}(\hat{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\hat{\mathbf{x}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}) \leq \mathcal{L}(\mathbf{x}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}).$$

The left saddle point inequality immediately implies

$$\sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\hat{\mathbf{x}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}) = f(\hat{\mathbf{x}}).$$

The right saddle point inequality, valid for a convex program under Slater's constraint qualification, entails

$$\mathcal{L}(\hat{\mathbf{x}}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \hat{\boldsymbol{\lambda}}, \hat{\boldsymbol{\mu}}) \leq \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

Hence, we can recover the minimum value of $f(\mathbf{x})$ as

$$f(\hat{\mathbf{x}}) = \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}).$$

The dual function

$$\mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

is jointly upper semicontinuous and concave in its arguments $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ and is well defined regardless of whether the program is convex. Maximization of the dual function subject to the constraint $\boldsymbol{\mu} \geq \mathbf{0}$ is referred to as the dual program. It trades an often simpler objective function in the original (primal) program for simpler constraints in the dual program.

In the absence of convexity or Slater's constraint qualification, we can still recover a weak form of duality based on the identity

$$f(\hat{\mathbf{x}}) = \inf_{\mathbf{x}} \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}),$$

which stems from the fact

$$\sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) = \begin{cases} f(\mathbf{x}) & \mathbf{x} \text{ feasible} \\ \infty & \mathbf{x} \text{ infeasible.} \end{cases}$$

Because $\inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq \mathcal{L}(\hat{\mathbf{x}}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\hat{\mathbf{x}})$, we can assert that

$$\sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \leq f(\hat{\mathbf{x}}) = \inf_{\mathbf{x}} \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}). \quad (15.4)$$

In other words, the minimum value of the primal problem exceeds the maximum value of the dual problem. Slater's constraint qualification guarantees for a convex program that the duality gap

$$\inf_{\mathbf{x}} \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu}) - \sup_{\boldsymbol{\lambda}, \boldsymbol{\mu} \geq \mathbf{0}} \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\mu})$$

vanishes when the primal problem has a finite minimum. Weak duality also makes it evident that if the primal program is unbounded below, then the dual program has no feasible point, and that if the dual program is unbounded above, then the primal program has no feasible point.

In a convex primal program, the equality constraint functions $g_i(\mathbf{x})$ are affine. If the inequality constraint functions $h_j(\mathbf{x})$ are also affine, then we can relate the dual program to the Fenchel conjugate $f^*(\mathbf{y})$ of $f(\mathbf{x})$. Suppose we write the constraints as $\mathbf{V}\mathbf{x} = \mathbf{d}$ and $\mathbf{W}\mathbf{x} \leq \mathbf{e}$. Then the dual function equals

$$\begin{aligned} \mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \inf_{\mathbf{x}} [f(\mathbf{x}) + \boldsymbol{\lambda}^*(\mathbf{V}\mathbf{x} - \mathbf{d}) + \boldsymbol{\mu}^*(\mathbf{W}\mathbf{x} - \mathbf{e})] \\ &= -\boldsymbol{\lambda}^*\mathbf{d} - \boldsymbol{\mu}^*\mathbf{e} + \inf_{\mathbf{x}} [f(\mathbf{x}) + (\mathbf{V}^*\boldsymbol{\lambda} + \mathbf{W}^*\boldsymbol{\mu})^*\mathbf{x}] \quad (15.5) \\ &= -\boldsymbol{\lambda}^*\mathbf{d} - \boldsymbol{\mu}^*\mathbf{e} - f^*(-\mathbf{V}^*\boldsymbol{\lambda} - \mathbf{W}^*\boldsymbol{\mu}). \end{aligned}$$

It may be that $f^*(\mathbf{y})$ equals ∞ for certain values of \mathbf{y} , but we can ignore these values in maximizing $\mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu})$. As pointed out in Proposition 14.3.1, $f^*(\mathbf{y})$ is a closed convex function.

The dual function may be differentiable even when the objective function or one of the constraints is not. Let $\boldsymbol{\omega} = (\boldsymbol{\lambda}, \boldsymbol{\mu})$, and suppose that the solution $\mathbf{x}(\boldsymbol{\omega})$ of $\mathcal{D}(\boldsymbol{\omega}) = \mathcal{L}(\mathbf{x}, \boldsymbol{\omega})$ is unique and depends continuously on $\boldsymbol{\omega}$. The inequalities

$$\begin{aligned} \mathcal{D}(\boldsymbol{\omega}_1) &= \mathcal{L}[\mathbf{x}(\boldsymbol{\omega}_1), \boldsymbol{\omega}_1] \leq \mathcal{L}[\mathbf{x}(\boldsymbol{\omega}_2), \boldsymbol{\omega}_1] \\ \mathcal{D}(\boldsymbol{\omega}_2) &= \mathcal{L}[\mathbf{x}(\boldsymbol{\omega}_2), \boldsymbol{\omega}_2] \leq \mathcal{L}[\mathbf{x}(\boldsymbol{\omega}_1), \boldsymbol{\omega}_2] \end{aligned}$$

can be re-expressed as

$$\begin{aligned} \sum_{i=1}^p [\lambda_{2i} - \lambda_{1i}]g_i[\mathbf{x}(\boldsymbol{\omega}_2)] + \sum_{j=1}^q [\mu_{2j} - \mu_{1j}]h_j[\mathbf{x}(\boldsymbol{\omega}_2)] &\leq \mathcal{D}(\boldsymbol{\omega}_2) - \mathcal{D}(\boldsymbol{\omega}_1) \\ \sum_{i=1}^p [\lambda_{2i} - \lambda_{1i}]g_i[\mathbf{x}(\boldsymbol{\omega}_1)] + \sum_{j=1}^q [\mu_{2j} - \mu_{1j}]h_j[\mathbf{x}(\boldsymbol{\omega}_1)] &\geq \mathcal{D}(\boldsymbol{\omega}_2) - \mathcal{D}(\boldsymbol{\omega}_1) \end{aligned}$$

Taking an appropriate convex combination of the left-hand sides of these two inequalities allows us to write $\mathcal{D}(\boldsymbol{\omega}_2) - \mathcal{D}(\boldsymbol{\omega}_1) = s(\boldsymbol{\omega}_2, \boldsymbol{\omega}_1)(\boldsymbol{\omega}_2 - \boldsymbol{\omega}_1)$ for a slope function $s(\boldsymbol{\omega}_2, \boldsymbol{\omega}_1)$. The slope requirement

$$\lim_{\boldsymbol{\omega}_2 \rightarrow \boldsymbol{\omega}_1} s(\boldsymbol{\omega}_2, \boldsymbol{\omega}_1) = [g_1(\mathbf{x}), \dots, g_p(\mathbf{x}), h_1(\mathbf{x}), \dots, h_q(\mathbf{x})]$$

with $\mathbf{x} = \mathbf{x}(\boldsymbol{\omega}_1)$ follows directly from the assumed continuity of the constraints and the solution vector $\mathbf{x}(\boldsymbol{\omega})$. Thus, $\mathcal{D}(\boldsymbol{\omega})$ is differentiable at $\boldsymbol{\omega}_1$.

15.5 Examples of Dual Programs

Here are some examples of dual programs that feature the close ties between duality and the Fenchel conjugate. We start with linear and quadratic programs, the simplest and most useful examples, and progress to more sophisticated examples.

Example 15.5.1 *Dual of a Linear Program*

The standard linear program minimizes $f(\mathbf{x}) = \mathbf{z}^* \mathbf{x}$ subject to the linear equality constraints $\mathbf{V} \mathbf{x} = \mathbf{d}$ and the nonnegativity constraints $\mathbf{x} \geq \mathbf{0}$. It is obvious that $f^*(\mathbf{y}) = \infty$ unless $\mathbf{y} = \mathbf{z}$, in which case $f^*(\mathbf{z}) = 0$. Thus in view of equation (15.5), the dual program maximizes $-\boldsymbol{\lambda}^* \mathbf{d}$ subject to the constraints $-\mathbf{V}^* \boldsymbol{\lambda} + \boldsymbol{\mu} = \mathbf{z}$ and $\boldsymbol{\mu} \geq \mathbf{0}$. Later we will demonstrate that there is no duality gap when either the primal or dual program has a finite solution [87, 183]. ■

Example 15.5.2 *Dual of a Strictly Convex Quadratic Program*

We have repeatedly visited the problem of minimizing the strictly convex function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ subject to the linear equality constraints $\mathbf{V} \mathbf{x} = \mathbf{d}$ and the linear inequality constraints $\mathbf{W} \mathbf{x} \leq \mathbf{e}$. Ignoring the constraints, $f(\mathbf{x})$ achieves its minimum value $-\frac{1}{2} \mathbf{b}^* \mathbf{A}^{-1} \mathbf{b} + c$ at the point $\mathbf{x} = -\mathbf{A}^{-1} \mathbf{b}$. In Example 14.3.1, we calculated the Fenchel conjugate

$$f^*(\mathbf{y}) = \frac{1}{2} (\mathbf{y} - \mathbf{b})^* \mathbf{A}^{-1} (\mathbf{y} - \mathbf{b}) - c.$$

If there are no equality constraints, then the dual program maximizes the quadratic

$$\begin{aligned} \mathcal{D}(\boldsymbol{\mu}) &= -\boldsymbol{\mu}^* \mathbf{e} - f^*(-\mathbf{W}^* \boldsymbol{\mu}) \\ &= -\boldsymbol{\mu}^* \mathbf{e} - \frac{1}{2} (\mathbf{W}^* \boldsymbol{\mu} + \mathbf{b})^* \mathbf{A}^{-1} (\mathbf{W}^* \boldsymbol{\mu} + \mathbf{b}) + c \\ &= -\frac{1}{2} \boldsymbol{\mu}^* \mathbf{W} \mathbf{A}^{-1} \mathbf{W}^* \boldsymbol{\mu} - (\mathbf{e} + \mathbf{W} \mathbf{A}^{-1} \mathbf{b})^* \boldsymbol{\mu} - \frac{1}{2} \mathbf{b}^* \mathbf{A}^{-1} \mathbf{b} + c \end{aligned}$$

subject to $\boldsymbol{\mu} \geq \mathbf{0}$. The dual program is easier to solve than the primal program when the number of inequality constraints is small and the number of variables is large.

In the presence of equality constraints and the absence of inequality constraints, the maximum of the dual function $\mathcal{D}(\boldsymbol{\lambda}) = -\boldsymbol{\lambda}^* \mathbf{d} - f^*(-\mathbf{V}^* \boldsymbol{\lambda})$ occurs where

$$\begin{aligned} \mathbf{0} &= -\mathbf{d} + \mathbf{V} \nabla f^*(-\mathbf{V}^* \boldsymbol{\lambda}) \\ &= -\mathbf{d} + \mathbf{V} \mathbf{A}^{-1} (-\mathbf{V}^* \boldsymbol{\lambda} - \mathbf{b}). \end{aligned}$$

If \mathbf{V} has full row rank, then the last equation has solution

$$\boldsymbol{\lambda} = -(\mathbf{V}\mathbf{A}^{-1}\mathbf{V}^*)^{-1}(\mathbf{V}\mathbf{A}^{-1}\mathbf{b} + \mathbf{d}).$$

This is just the Lagrange multiplier calculated in Example 5.2.6 and Proposition 5.2.2. Given the multiplier, it is straightforward to calculate the optimal value of the primal variable \mathbf{x} . ■

Example 15.5.3 *Dual of a Linear Semidefinite Program*

The linear semidefinite programming problem consists in minimizing the trace function $\mathbf{X} \mapsto \text{tr}(\mathbf{C}\mathbf{X})$ over the cone of positive semidefinite matrices S_+^n subject to the linear constraints $\text{tr}(\mathbf{A}_i\mathbf{X}) = b_i$ for $1 \leq i \leq p$. Here \mathbf{C} and the \mathbf{A}_i are assumed symmetric. According to Sylvester’s criterion, the constraint $\mathbf{B} \in S_+^n$ involves a complicated system of nonlinear inequalities. It is conceptually simpler to rewrite the constraint $\mathbf{B} \in S_+^n$ as $-\lambda_1(\mathbf{X}) \leq 0$, where $\lambda_1(\mathbf{X})$ is the minimum eigenvalue of \mathbf{X} . Example 6.3.8 proves that $\lambda_1(\mathbf{X})$ is concave in \mathbf{X} .

In defining the dual problem, there is a generic way of accommodating complicated convex constraints. Suppose in the standard convex program we confine \mathbf{x} to some closed convex set S in addition to imposing explicit equality and inequality constraints. Minimizing the objective function $f(\mathbf{x})$ subject to all of the constraints is equivalent to minimizing $f(\mathbf{x}) + \delta_S(\mathbf{x})$ subject to the functional constraints alone. In forming the dual we therefore take the infimum of the Lagrangian over \mathbf{x} in S rather than over all \mathbf{x} .

In linear semidefinite programming, we therefore define the dual function

$$\begin{aligned} \mathcal{D}(\boldsymbol{\lambda}) &= \inf_{\mathbf{X} \in S_+^n} \left\{ \text{tr}(\mathbf{C}\mathbf{X}) + \sum_{i=1}^p \lambda_i [b_i - \text{tr}(\mathbf{A}_i\mathbf{X})] \right\} \\ &= \mathbf{b}^* \boldsymbol{\lambda} + \inf_{\mathbf{X} \in S_+^n} \text{tr} \left[\left(\mathbf{C} - \sum_{i=1}^p \lambda_i \mathbf{A}_i \right) \mathbf{X} \right] \\ &= \mathbf{b}^* \boldsymbol{\lambda} + \inf_{\mathbf{X} \in S_+^n} t(\mathbf{X}). \end{aligned}$$

If $t(\mathbf{X}) < 0$ for some $\mathbf{X} \in S_+^n$, then $t(\mathbf{X})$ can be made to approach $-\infty$ by replacing \mathbf{X} by $c\mathbf{X}$ and taking c large. It follows that we should restrict the matrix $\mathbf{C} - \sum_{i=1}^p \lambda_i \mathbf{A}_i$ to lie in the negative of the polar cone $(S_+^n)^\circ$ of S_+^n . We know from Example 14.3.7 that $(S_+^n)^\circ = -S_+^n$. When this restriction holds and $\mathbf{C} - \sum_{i=1}^p \lambda_i \mathbf{A}_i$ is positive semidefinite, the minimum of $t(\mathbf{X})$ is achieved by setting $\mathbf{X} = \mathbf{0}$. The dual problem thus becomes one of maximizing $\mathbf{b}^* \boldsymbol{\lambda}$ subject to the condition that $\mathbf{C} - \sum_{i=1}^p \lambda_i \mathbf{A}_i$ is positive semidefinite. Slater’s constraint qualification guaranteeing a duality gap of 0 is just the assumption that there exists a positive definite matrix \mathbf{X} that is feasible for the primal problem. ■

Example 15.5.4 *Regression with a Non-Euclidean Norm*

In unconstrained least squares one minimizes the criterion $\|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|$ with respect to $\boldsymbol{\beta}$. If we substitute a non-Euclidean norm $\|\mathbf{z}\|_{\dagger}$ for the Euclidean norm $\|\mathbf{z}\|$, then the problem becomes much harder. Suppose the dual norm to $\|\mathbf{z}\|_{\dagger}$ is $\|\mathbf{w}\|_{\star}$. With no Lagrange multipliers in sight, the dual function is constant. However, let us repose the problem as minimizing $\|\mathbf{z}\|_{\dagger}$ subject to the linear constraint $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$. The Lagrangian for this version of the problem is obviously $\mathcal{L}(\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\lambda}) = \|\mathbf{z}\|_{\dagger} + \boldsymbol{\lambda}^*(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z})$. In view of Example 14.3.5, the Fenchel conjugate of $f(\mathbf{z}, \boldsymbol{\beta}) = \|\mathbf{z}\|_{\dagger}$ is

$$f^*(\mathbf{w}, \gamma) = \sup_{\mathbf{z}, \boldsymbol{\beta}} (\mathbf{w}^* \mathbf{z} + \gamma^* \boldsymbol{\beta} - \|\mathbf{z}\|_{\dagger}) = \delta_B(\mathbf{w}) + \delta_{\{0\}}(\gamma),$$

where B is the closed unit ball associated with the dual norm. Equation (15.5) therefore produces the dual function

$$\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^* \mathbf{y} - \delta_B(\boldsymbol{\lambda}) - \delta_{\{0\}}(\mathbf{X}^* \boldsymbol{\lambda}).$$

The dual problem consists of minimizing $-\boldsymbol{\lambda}^* \mathbf{y}$ subject to the constraints $\mathbf{X}^* \boldsymbol{\lambda} = \mathbf{0}$ and $\|\boldsymbol{\lambda}\|_{\star} \leq 1$.

We can reformulate the regression problem as minimizing the criterion $\frac{1}{2} \|\mathbf{z}\|_{\dagger}^2$ subject to the linear constraint $\mathbf{z} = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}$. The Lagrangian now becomes

$$\mathcal{L}(\mathbf{z}, \boldsymbol{\beta}, \boldsymbol{\lambda}) = \frac{1}{2} \|\mathbf{z}\|_{\dagger}^2 + \boldsymbol{\lambda}^*(\mathbf{y} - \mathbf{X}\boldsymbol{\beta} - \mathbf{z}),$$

and in view of Example 14.3.5 the dual function reduces to

$$\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^* \mathbf{y} - \frac{1}{2} \|\boldsymbol{\lambda}\|_{\star}^2 - \delta_{\{0\}}(\mathbf{X}^* \boldsymbol{\lambda}).$$

Hence, the dual problem consists of minimizing the function $-\boldsymbol{\lambda}^* \mathbf{y} + \frac{1}{2} \|\boldsymbol{\lambda}\|_{\star}^2$ subject to the single constraint $\mathbf{X}^* \boldsymbol{\lambda} = \mathbf{0}$. In this case there are two different dual problems. ■

Example 15.5.5 *Dual of a Geometric Program*

In passing to the dual, it is helpful to restate geometric programming as

$$\begin{aligned} &\text{minimize} && \ln \left(\sum_{k=1}^{m_0} e^{\mathbf{a}_{0k}^* \mathbf{x} + b_{0k}} \right) \\ &\text{subject to} && \ln \left(\sum_{k=1}^{m_j} e^{\mathbf{a}_{jk}^* \mathbf{x} + b_{jk}} \right) \leq 0, \quad 1 \leq j \leq q. \end{aligned}$$

In this convex version of geometric programming, the positive constants multiplying the monomials are just the exponentials $e^{b_{jk}}$. It is hard to take the Fenchel conjugate of the associated Lagrangian so we resort to the standard trick of introducing simpler functions $f_j(\mathbf{z}_j)$ with distinct arguments

and compensating equality constraints. If $f_j(\mathbf{z}_j)$ equals $\ln(\sum_k e^{z_{jk}})$, then Example 14.3.2 calculates the entropy conjugate

$$f_j^*(\mathbf{y}_j) = \begin{cases} \sum_k y_{jk} \ln y_{jk} & \text{for } \sum_k y_{jk} = 1 \text{ and all } y_{jk} \geq 0 \\ \infty & \text{otherwise.} \end{cases}$$

The equality constraint corresponding to $f_j(\mathbf{z}_j)$ is $\mathbf{A}_j \mathbf{x} + \mathbf{b}_j - \mathbf{z}_j = \mathbf{0}$ for a matrix \mathbf{A}_j with rows \mathbf{a}_{jk}^* . In this notation, the simplified Lagrangian becomes

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \mathbf{z}_0, \dots, \mathbf{z}_q, \boldsymbol{\lambda}, \boldsymbol{\mu}) &= f_0(\mathbf{z}_0) + \sum_{j=1}^q \mu_j f_j(\mathbf{z}_j) + \sum_{j=0}^q \boldsymbol{\lambda}_j^* (\mathbf{A}_j \mathbf{x} + \mathbf{b}_j - \mathbf{z}_j) \\ &= f_0(\mathbf{z}_0) - \boldsymbol{\lambda}_0^* \mathbf{z}_0 + \sum_{j=1}^q \mu_j [f_j(\mathbf{z}_j) - \mu_j^{-1} \boldsymbol{\lambda}_j^* \mathbf{z}_j] + \sum_{j=0}^q \boldsymbol{\lambda}_j^* (\mathbf{A}_j \mathbf{x} + \mathbf{b}_j). \end{aligned}$$

It follows that the dual function equals

$$\begin{aligned} \mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \inf_{\mathbf{x}, \mathbf{z}_0, \dots, \mathbf{z}_q} \mathcal{L}(\mathbf{x}, \mathbf{z}_0, \dots, \mathbf{z}_q, \boldsymbol{\lambda}, \boldsymbol{\mu}) \\ &= \sum_{j=0}^q \boldsymbol{\lambda}_j^* \mathbf{b}_j - f_0^*(\boldsymbol{\lambda}_0) - \sum_{j=1}^q \mu_j f_j^*(\mu_j^{-1} \boldsymbol{\lambda}_j) + \delta_{\mathbf{0}} \left(\sum_{j=0}^q \mathbf{A}_j^* \boldsymbol{\lambda}_j \right). \end{aligned}$$

The exceptional cases where one or more $\mu_j = 0$ require special treatment. When $\mu_j = 0$ but $\boldsymbol{\lambda}_j \neq \mathbf{0}$, the dual function must be interpreted as $-\infty$. When $\mu_j = 0$ and $\boldsymbol{\lambda}_j = \mathbf{0}$, the expression for $\mathcal{D}(\boldsymbol{\lambda}, \boldsymbol{\mu})$ is valid provided we interpret $\mu_j f_j^*(\mu_j^{-1} \boldsymbol{\lambda}_j)$ as 0. Because of the nature of the entropy function, the constraints $\boldsymbol{\lambda}_j \geq \mathbf{0}$ for $0 \leq j \leq q$ and $\|\boldsymbol{\lambda}_0\|_1 = 1$ and $\|\boldsymbol{\lambda}_j\|_1 = \mu_j$ for $1 \leq j \leq q$ are implicit in the dual function. The constraints $\mu_j \geq 0$ for $1 \leq j \leq q$ are explicit, as is the constraint $\sum_{j=0}^q \mathbf{A}_j^* \boldsymbol{\lambda}_j = \mathbf{0}$. ■

Example 15.5.6 *Dijkstra’s Algorithm as Block Relaxation of the Dual*

Dijkstra’s problem can be restated as finding the minimum of the convex function

$$g(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=0}^{r-1} \delta_{C_i}(\mathbf{x}) = f(\mathbf{x}) + \delta_C(\mathbf{x})$$

for $f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|^2$, C the intersection of the closed convex sets C_0 through C_{r-1} , $\delta_C(\mathbf{x})$ the indicator function of C as defined in Example 14.3.4, and \mathbf{z} the external point to C . As in Example 15.5.5, it is difficult to take the Fenchel conjugate of $g(\mathbf{x})$. Matters simplify tremendously if we replace the argument of each $\delta_{C_i}(\mathbf{x})$ by \mathbf{x}_i and impose the constraint $\mathbf{x}_i = \mathbf{x}$. Consider therefore the Lagrangian

$$\mathcal{L}(\mathbf{X}, \boldsymbol{\Lambda}) = f(\mathbf{x}) + \sum_{i=0}^{r-1} \delta_{C_i}(\mathbf{x}_i) + \sum_{i=0}^{r-1} \boldsymbol{\lambda}_i^* (\mathbf{x} - \mathbf{x}_i)$$

with $\mathbf{X} = (\mathbf{x}, \mathbf{x}_0, \dots, \mathbf{x}_{r-1})$ and $\mathbf{\Lambda} = (\boldsymbol{\lambda}_0, \dots, \boldsymbol{\lambda}_{r-1})$. The dual function is

$$\begin{aligned} \mathcal{D}(\mathbf{\Lambda}) &= -\sup_{\mathbf{X}} \left[-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i^* \mathbf{x} - f(\mathbf{x}) + \sum_{i=0}^{r-1} \boldsymbol{\lambda}_i^* \mathbf{x}_i - \sum_{i=0}^{r-1} \delta_{C_i}(\mathbf{x}_i) \right] \\ &= -f^* \left(-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i \right) - \sum_{i=0}^{r-1} \delta_{C_i}^*(\boldsymbol{\lambda}_i). \end{aligned}$$

The dual problem consists of minimizing the convex function

$$f^* \left(-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i \right) + \sum_{i=0}^{r-1} \delta_{C_i}^*(\boldsymbol{\lambda}_i).$$

Dykstra's algorithm solves the dual problem by block descent [9, 26].

Suppose that we fix all $\boldsymbol{\lambda}_i$ except $\boldsymbol{\lambda}_j$. The stationarity condition requires $\mathbf{0}$ to belong to the subdifferential

$$\partial \left[f^* \left(-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i \right) + \delta_{C_j}^*(\boldsymbol{\lambda}_j) \right] = -\partial f^* \left(-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i \right) + \partial \delta_{C_j}^*(\boldsymbol{\lambda}_j).$$

It follows that there exists a vector \mathbf{x}_j such that $\mathbf{x}_j \in \partial f^* \left(-\sum_{i=0}^{r-1} \boldsymbol{\lambda}_i \right)$ and $\mathbf{x}_j \in \partial \delta_{C_j}^*(\boldsymbol{\lambda}_j)$. Propositions 14.3.1 and 14.4.4 allow us to invert these two relations. Thus, $-\boldsymbol{\lambda}_j \in \partial [f(\mathbf{x}_j) + \sum_{i \neq j} \boldsymbol{\lambda}_i^* \mathbf{x}_j]$ and $\boldsymbol{\lambda}_j \in \partial \delta_{C_j}(\mathbf{x}_j)$, which together are equivalent to the primal stationarity condition

$$\mathbf{0} \in \partial \left[f(\mathbf{x}_j) + \sum_{i \neq j} \boldsymbol{\lambda}_i^* \mathbf{x}_j + \delta_{C_j}(\mathbf{x}_j) \right].$$

As a consequence, it suffices to minimize

$$f(\mathbf{x}_j) + \sum_{i \neq j} \boldsymbol{\lambda}_i^* \mathbf{x}_j + \delta_{C_j}(\mathbf{x}_j) = \frac{1}{2} \left\| \mathbf{x}_j + \sum_{i \neq j} \boldsymbol{\lambda}_i - \mathbf{z} \right\|^2 + \delta_{C_j}(\mathbf{x}_j) + c,$$

where c is an irrelevant constant. But this problem is solved by projecting $\mathbf{z} - \sum_{i \neq j} \boldsymbol{\lambda}_i$ onto the convex set C_j .

The update of $\boldsymbol{\lambda}_j$ satisfies

$$\boldsymbol{\lambda}_j = -\partial \left[f(\mathbf{x}_j) + \sum_{i \neq j} \boldsymbol{\lambda}_i^* \mathbf{x}_j \right] = \mathbf{z} - \mathbf{x}_j - \sum_{i \neq j} \boldsymbol{\lambda}_i.$$

Given the converged values of the $\boldsymbol{\lambda}_j$, the optimal \mathbf{x} can be recovered from the stationarity condition

$$\mathbf{0} = \partial f(\mathbf{x}) + \sum_{i=0}^{r-1} \boldsymbol{\lambda}_i = \mathbf{x} - \mathbf{z} + \sum_{i=0}^{r-1} \boldsymbol{\lambda}_i$$

for the Lagrangian $\mathcal{L}(\mathbf{X}, \mathbf{\Lambda})$ as $\mathbf{x} = \mathbf{z} - \sum_{i=0}^{r-1} \boldsymbol{\lambda}_i$ ■

Example 15.5.7 *Duffin's Counterexample*

Consider the convex program of minimizing $f(\mathbf{x}) = e^{-x_2}$ subject to the inequality constraint $h(\mathbf{x}) = \|\mathbf{x}\| - x_1 \leq 0$ on \mathbb{R}^2 . This problem does not satisfy Slater's condition because all feasible \mathbf{x} satisfy $x_2 = 0$ and consequently $h(\mathbf{x}) = 0$ and $f(\mathbf{x}) = 1$. To demonstrate that there is a duality gap, we show that the dual function

$$\mathcal{D}(\mu) = \inf_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mu) = \inf_{\mathbf{x}} [e^{-x_2} + \mu(\|\mathbf{x}\| - x_1)]$$

is identically 0. Because $\mathcal{L}(\mathbf{x}, \mu) \geq 0$ for all \mathbf{x} and $\mu \geq 0$, it suffices to prove that $\mathcal{L}(\mathbf{x}, \mu)$ can be made less than any positive ϵ . Choose an x_2 so that $e^{-x_2} < \epsilon/2$. Having chosen x_2 , choose x_1 so that

$$\begin{aligned} \sqrt{x_1^2 + x_2^2} - x_1 &= x_1 \sqrt{1 + \frac{x_2^2}{x_1^2}} - x_1 \\ &\leq x_1 \left(1 + \frac{x_2^2}{2x_1^2}\right) - x_1 \\ &= \frac{x_2^2}{2x_1} \\ &< \frac{\epsilon}{2\mu}. \end{aligned}$$

With these choices we have $\mathcal{L}(\mathbf{x}, \mu) < \epsilon$. Thus, the minimum value 1 of the primal problem is strictly greater than the maximum value 0 of the dual problem. ■

Before turning to practical applications of duality in the next section, we would like to mention the fundamental theorem of linear programming. Our proof depends on the subtle properties of polyhedral sets sketched in Appendix A.3. Readers can digest this material at their leisure.

Proposition 15.5.1 *If either the primal or dual linear program formulated in Example 15.5.1 has a solution, then the other program has a solution as well. Furthermore, there is no duality gap.*

Proof: If the primal program has a solution, then inequality (15.4) shows that the dual program has an upper bound. Proposition A.3.5 therefore implies that the dual program has a solution. Conversely, if the dual program has a solution, then inequality (15.4) shows that the primal program has a lower bound. Since a linear function is both convex and concave, a second application of Proposition A.3.5 shows that the primal program has a solution. According to Example 6.5.5, the existence of either solution forces the preferred form of the Lagrange multiplier rule and hence implies no duality gap. ■

Example 15.5.8 *Von Neumann's Minimax Theorem*

Von Neumann's minimax theorem is one of the earliest results of game theory [266]. In purely mathematical terms, it can be stated as the identity

$$\min_{\mathbf{y} \in S} \max_{\mathbf{x} \in S} \mathbf{x}^* \mathbf{A} \mathbf{y} = \max_{\mathbf{x} \in S} \min_{\mathbf{y} \in S} \mathbf{x}^* \mathbf{A} \mathbf{y}, \tag{15.6}$$

where $\mathbf{A} = (a_{ij})$ is an $n \times n$ matrix and S is the unit simplex

$$S = \left\{ \mathbf{z} \in \mathbb{R}^n : z_1 \geq 0, \dots, z_n \geq 0, \sum_{i=1}^n z_i = 1 \right\}.$$

It is possible to view the identity (15.6) as a manifestation of linear programming duality. As the reader can check (Problem 28), the primal program

$$\begin{aligned} & \text{minimize} && u \\ & \text{subject to} && \sum_{j=1}^n y_j = 1, \sum_{j=1}^n a_{ij} y_j \leq u \quad \forall i, y_j \geq 0 \quad \forall j \end{aligned}$$

has dual program

$$\begin{aligned} & \text{maximize} && v \\ & \text{subject to} && \sum_{i=1}^n x_i = 1, \sum_{i=1}^n x_i a_{ij} \geq v \quad \forall j, x_i \geq 0 \quad \forall i. \end{aligned}$$

Von Neumann's identity (15.6) is true because the primal and optimal values p and d of this linear program satisfy

$$\begin{aligned} \min_{\mathbf{y} \in S} \max_{\mathbf{x} \in S} \mathbf{x}^* \mathbf{A} \mathbf{y} &= \min_{\mathbf{y} \in S} \max_{1 \leq i \leq n} \mathbf{e}_i^* \mathbf{A} \mathbf{y} = p \\ \max_{\mathbf{x} \in S} \min_{\mathbf{y} \in S} \mathbf{x}^* \mathbf{A} \mathbf{y} &= \max_{\mathbf{x} \in S} \min_{1 \leq j \leq n} \mathbf{x}^* \mathbf{A} \mathbf{e}_j = d. \end{aligned}$$

Von Neumann's identity is a special case of the much more general minimax principle of Sion [155, 238]. This principle implies no duality gap in convex programming given appropriate compactness assumptions. ■

15.6 Practical Applications of Duality

The two examples of this section illustrate how duality is not just a theoretical construct. It can also lead to the discovery of concrete optimization algorithms. Practitioners of optimization should always bear this in mind as well as the lower bound offered by inequality (15.4). Of course, solving the dual problem is seldom enough. To realize the potential of duality, one must convert the solution of the dual problem into a solution of the primal problem.

Example 15.6.1 *The Power Plant Problem*

The power plant production problem [226] involves minimizing

$$f(\mathbf{x}) = \sum_{i=1}^n f_i(x_i)$$

subject to the constraints $0 \leq x_i \leq u_i$ for each i and $\sum_{i=1}^n x_i \geq d$. For plant i , x_i is the power output, u_i is the capacity, and $f_i(x_i)$ is the cost. The total demand is d . The Lagrangian for this minimization problem is

$$\mathcal{L}(\mathbf{x}, \mu) = \sum_{i=1}^n f_i(x_i) + \mu \left(d - \sum_{i=1}^n x_i \right).$$

As a consequence of the separability of the Lagrangian, the dual function can be expressed as

$$\mathcal{D}(\mu) = \mu d + \sum_{i=1}^n \min_{0 \leq x_i \leq u_i} [f_i(x_i) - \mu x_i].$$

For the quadratic choices $f_i(x_i) = a_i x_i + \frac{1}{2} b_i x_i^2$ with positive cost constants a_i and b_i , the problem is a convex program, and it is possible to explicitly solve for the dual. A brief calculation shows that optimal value of x_i is

$$\hat{x}_i = \begin{cases} 0 & 0 \leq \mu \leq a_i \\ \frac{\mu - a_i}{b_i} & a_i \leq \mu \leq a_i + b_i u_i \\ u_i & \mu \geq a_i + b_i u_i. \end{cases}$$

These solutions translate into the dual function

$$\mathcal{D}(\mu) = \mu d + \sum_{i=1}^n \begin{cases} 0 & 0 \leq \mu \leq a_i \\ -\frac{(\mu - a_i)^2}{2b_i} & a_i \leq \mu \leq a_i + b_i u_i \\ a_i u_i + \frac{1}{2} b_i u_i^2 - \mu u_i & \mu \geq a_i + b_i u_i. \end{cases}$$

and ultimately into the derivative $\mathcal{D}'(\mu) = d - \sum_{i=1}^n \hat{x}_i$. Because $\mathcal{D}(\mu)$ is concave, a stationary point furnishes the global maximum. It is straightforward to implement bisection to locate a stationary point. Steepest ascent is another route to maximizing $\mathcal{D}(\mu)$. Problem 26 asks the reader to consider the impact of assuming one or more $b_i = 0$. ■

Example 15.6.2 *Linear Classification*

Classification problems are ubiquitous in statistics. Section 13.8 discusses one approach to discriminant analysis. Here we take another motivated by hyperplane separation. The binary classification problem can be phrased in terms of a training sequence of observation vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ from \mathbb{R}^n and

an associated sequence of population indicators s_1, \dots, s_m from $\{-1, +1\}$. In favorable situations, the two different populations can be separated by a hyperplane defined by a unit vector \mathbf{z} and constants $c_1 \leq c_2$ in the sense that

$$\begin{aligned}\mathbf{z}^* \mathbf{v}_i &\leq c_1, & s_i &= -1 \\ \mathbf{z}^* \mathbf{v}_i &\geq c_2, & s_i &= +1.\end{aligned}$$

The optimal separation occurs when the difference $c_2 - c_1$ is maximized.

This linear classification problem can be simplified by rewriting the separation conditions as

$$\begin{aligned}\mathbf{z}^* \mathbf{v}_i - \frac{c_1 + c_2}{2} &\leq -\frac{c_2 - c_1}{2}, & s_i &= -1 \\ \mathbf{z}^* \mathbf{v}_i - \frac{c_1 + c_2}{2} &\geq +\frac{c_2 - c_1}{2}, & s_i &= +1.\end{aligned}$$

If we let $a = (c_2 - c_1)/2$, $b = (c_1 + c_2)/(c_2 - c_1)$, and $\mathbf{y} = a^{-1}\mathbf{z}$, then these become

$$\begin{aligned}\mathbf{y}^* \mathbf{v}_i - b &\leq -1, & s_i &= -1 \\ \mathbf{y}^* \mathbf{v}_i - b &\geq +1, & s_i &= +1.\end{aligned}\tag{15.7}$$

Thus, the linear classification problem reduces to minimizing the criterion $\frac{1}{2}\|\mathbf{y}\|^2$ subject to the inequality constraints (15.7). Observe that the constraint functions are linear in the parameter vector $(\mathbf{y}^*, b)^*$ in this semidefinite quadratic programming problem. Unfortunately, because the component b does not appear in the objective function $\frac{1}{2}\|\mathbf{y}\|^2$, Dykstra's algorithm does not apply. Once we find \mathbf{y} and b , we can classify a new test vector \mathbf{v} in the $s = -1$ population when $\mathbf{y}^* \mathbf{v} - b < 0$ and in the $s = +1$ population when $\mathbf{y}^* \mathbf{v} - b > 0$.

There is no guarantee that a feasible vector $(\mathbf{y}^*, b)^*$ exists for the linear classification problem as stated. A more realistic version of the problem imposes the inequality constraints

$$s_i(\mathbf{y}^* \mathbf{v}_i - b) \geq 1 - \epsilon_i\tag{15.8}$$

using a slack variable $\epsilon_i \geq 0$. To penalize deviation from the ideal of perfect separation by a hyperplane, we modify the objective function to be

$$f(\mathbf{y}, b, \epsilon) = \frac{1}{2}\|\mathbf{y}\|^2 + \delta \sum_{i=1}^m \epsilon_i\tag{15.9}$$

for some tuning constant $\delta > 0$. The constraints (15.8) and $\epsilon_i \geq 0$ are again linear in the parameter vector $\mathbf{x} = (\mathbf{y}^*, b, \epsilon^*)^*$. The Lagrangian

$$\begin{aligned} \mathcal{L}(\mathbf{y}, b, \boldsymbol{\epsilon}, \boldsymbol{\mu}) &= \frac{1}{2} \|\mathbf{y}\|^2 + \delta \sum_{i=1}^m \epsilon_i - \sum_{i=1}^m \mu_{m+i} \epsilon_i \\ &\quad + \sum_{i=1}^m \mu_i [-s_i (\mathbf{v}_i^* \mathbf{y} - b) + 1 - \epsilon_i] \end{aligned}$$

involves $2m$ nonnegative multipliers μ_1, \dots, μ_{2m} .

It is simpler to solve the dual problem than the primal problem. We can formulate the dual by following the steps of Example 15.5.2. If we express the inequality constraints as $\mathbf{W}\mathbf{x} \leq \mathbf{e}$, then the matrix transpose \mathbf{W}^* and the vector \mathbf{e} are

$$\mathbf{W}^* = - \begin{pmatrix} s_1 \mathbf{v}_1 & \cdots & s_m \mathbf{v}_m & \mathbf{0} \\ -s_1 & \cdots & -s_m & 0 \\ & & \mathbf{I}_m & \mathbf{I}_m \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} -1 \\ \mathbf{0} \end{pmatrix}$$

in the current setting. The dual problem consists of maximizing the function $-\boldsymbol{\mu}^* \mathbf{e} - f^*(-\mathbf{W}^* \boldsymbol{\mu})$ subject to the constraints $\mu_i \geq 0$ and restrictions imposed by the essential domain of the Fenchel conjugate $f^*(\mathbf{p}, q, \mathbf{r})$. An easy calculation gives

$$\begin{aligned} f^*(\mathbf{p}, q, \mathbf{r}) &= \sup_{\mathbf{y}, b, \boldsymbol{\epsilon}} \left[\mathbf{p}^* \mathbf{y} + qb + \mathbf{r}^* \boldsymbol{\epsilon} - \frac{1}{2} \|\mathbf{y}\|^2 - \delta \sum_{i=1}^m \epsilon_i \right] \\ &= \begin{cases} \infty & q \neq 0 \text{ or } r_i \neq \delta \text{ for some } i \\ \mathbf{p}^* \mathbf{p} - \frac{1}{2} \|\mathbf{p}\|^2 & \text{otherwise} \end{cases} \\ &= \begin{cases} \infty & q \neq 0 \text{ or } r_i \neq \delta \text{ for some } i \\ \frac{1}{2} \|\mathbf{p}\|^2 & \text{otherwise.} \end{cases} \end{aligned}$$

To match $-\mathbf{W}^* \boldsymbol{\mu}$ to the indicated essential domain, note that the restriction $q = 0$ entails the constraint $\sum_{i=1}^m s_i \mu_i = 0$, and the restriction $r_i = \delta$ entails the constraint $\mu_i + \mu_{m+i} = \delta$ and therefore the bound $\mu_i \leq \delta$. For the vector \mathbf{p} we substitute the linear combination $\sum_{i=1}^m s_i \mu_i \mathbf{v}_i$. Hence, the dual problem consists in maximizing

$$\begin{aligned} -\boldsymbol{\mu}^* \mathbf{e} - f^*(-\mathbf{W}^* \boldsymbol{\mu}) &= \sum_{i=1}^m \mu_i - \frac{1}{2} \left\| \sum_{i=1}^m s_i \mu_i \mathbf{v}_i \right\|^2 \tag{15.10} \\ &= \sum_{i=1}^m \mu_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m s_i \mu_i \mathbf{v}_i^* \mathbf{v}_j s_j \mu_j \end{aligned}$$

subject to the constraints $\sum_{i=1}^m s_i \mu_i = 0$ and $0 \leq \mu_i \leq \delta$ for all i .

Solving the dual problem fortunately leads to a straightforward solution of the primal problem. For instance, the Lagrangian conditions

$$\frac{\partial}{\partial y_j} \mathcal{L}(\mathbf{y}, b, \boldsymbol{\epsilon}, \boldsymbol{\mu}) = 0$$

give $\mathbf{y} = \sum_{i=1}^m \mu_i s_i \mathbf{v}_i$. The Lagrangian condition

$$\frac{\partial}{\partial \epsilon_j} \mathcal{L}(\mathbf{y}, b, \epsilon, \boldsymbol{\mu}) = \delta - \mu_j - \mu_{m+j} = 0$$

implies that $\mu_{m+j} = \delta - \mu_j$. The complementary slackness conditions

$$\begin{aligned} 0 &= \mu_i [-s_i(\mathbf{v}_i^* \mathbf{y} - b) + 1 - \epsilon_i] \\ 0 &= \mu_{m+j} \epsilon_j \end{aligned}$$

can be used to determine b and ϵ_i for $1 \leq i \leq m$. If we choose an index j such that $0 < \mu_j < \delta$, then $\mu_{m+j} > 0$ and $\epsilon_j = 0$. It follows that $-s_j(\mathbf{v}_j^* \mathbf{y} - b) + 1 = 0$ and that $b = \mathbf{v}_j^* \mathbf{y} - s_j$ since $s_j = \pm 1$. Given b , all ϵ_j with $\mu_j > 0$ are determined. If $\mu_j = 0$, then $\mu_{m+j} = \delta > 0$ and $\epsilon_j = 0$.

Despite these interesting maneuvers, we have not actually shown how to solve the dual problem. For linear classification problems with many training vectors $\mathbf{v}_1, \dots, \mathbf{v}_m$ in a high-dimensional space \mathbb{R}^n , it is imperative to formulate an efficient algorithm. One possibility is to try a penalty method. If we subtract the square penalty $\omega(\sum_{i=1}^m s_i \mu_i)^2$ from the dual function (15.10), then the remaining box constraints are consistent with coordinate ascent. Sending ω to ∞ then produces the dual solution. Alternatively, subtracting the penalty $\omega|\sum_{i=1}^m s_i \mu_i|$ converts the dual problem into an exact penalty problem and opens up the possibility of path following as described in Chap. 16. ■

15.7 Problems

1. Let $C \subset \mathbb{R}^2$ be the cone defined by the constraint $x_1 \leq x_2$. Show that the projection operator $P_C(\mathbf{x})$ has components

$$P_C(\mathbf{x})_j = \begin{cases} x_j & x_1 \leq x_2 \\ \frac{1}{2}(x_1 + x_2) & x_1 > x_2. \end{cases}$$

Let $S \subset \mathbb{R}^3$ be the cone defined by the constraint $x_2 \leq \frac{1}{2}(x_1 + x_3)$. Show that the projection operator $P_S(\mathbf{x})$ has components

$$P_S(\mathbf{x})_j = \begin{cases} x_j & x_2 \leq \frac{1}{2}(x_1 + x_3) \\ \frac{5}{6}x_1 + \frac{1}{3}x_2 - \frac{1}{6}x_3 & j = 1 \text{ and } x_2 > \frac{1}{2}(x_1 + x_3) \\ \frac{1}{3}x_1 + \frac{1}{3}x_2 + \frac{1}{3}x_3 & j = 2 \text{ and } x_2 > \frac{1}{2}(x_1 + x_3) \\ -\frac{1}{6}x_1 + \frac{1}{3}x_2 + \frac{5}{6}x_3 & j = 3 \text{ and } x_2 > \frac{1}{2}(x_1 + x_3). \end{cases}$$

2. If A and B are two closed convex sets, then prove that projection onto the Cartesian product $A \times B$ is effected by the Cartesian product operator $(\mathbf{x}, \mathbf{y}) \mapsto [P_A(\mathbf{x}), P_B(\mathbf{y})]$.

3. Program and test either an algorithm for projection onto the unit simplex or the pool adjacent violators algorithm.
4. If C is a closed convex set in \mathbb{R}^n and $\mathbf{x} \notin C$, then demonstrate that

$$\text{dist}(\mathbf{x}, C) = \inf_{\mathbf{y} \in C} \sup_{\|\mathbf{z}\|=1} \mathbf{z}^*(\mathbf{x} - \mathbf{y}) = \sup_{\|\mathbf{z}\|=1} \inf_{\mathbf{y} \in C} \mathbf{z}^*(\mathbf{x} - \mathbf{y}).$$

Also prove that there exists a unit vector \mathbf{z} with

$$\text{dist}(\mathbf{x}, C) = \inf_{\mathbf{y} \in C} \mathbf{z}^*(\mathbf{x} - \mathbf{y}).$$

(Hints: The first equality follows from the Cauchy-Schwarz inequality and the definition of $\text{dist}(\mathbf{x}, C)$. The rest of the problem depends on the Cauchy-Schwarz inequality, the particular choice

$$\mathbf{z} = \text{dist}(\mathbf{x}, C)^{-1}[\mathbf{x} - P_C(\mathbf{x})],$$

and the obtuse angle criterion.)

5. Let S and T be subspaces of \mathbb{R}^n . Demonstrate that the projections P_S and P_T satisfy $P_S P_T = P_{S \cap T}$ if and only if $P_S P_T = P_T P_S$.
6. Let C be a closed convex set in \mathbb{R}^n . Show that
- $\text{dist}(\mathbf{x} + \mathbf{y}, C + \mathbf{y}) = \text{dist}(\mathbf{x}, C)$ for all \mathbf{x} and \mathbf{y} .
 - $P_{C+\mathbf{y}}(\mathbf{x} + \mathbf{y}) = P_C(\mathbf{x}) + \mathbf{y}$ for all \mathbf{x} and \mathbf{y} .
 - $\text{dist}(a\mathbf{x}, aC) = |a| \text{dist}(\mathbf{x}, C)$ for all \mathbf{x} and real a .
 - $P_{aC}(a\mathbf{x}) = aP_C(\mathbf{x})$ for all \mathbf{x} and real a .

Let S be a subspace of \mathbb{R}^n . Show that

- $\text{dist}(\mathbf{x} + \mathbf{y}, S) = \text{dist}(\mathbf{x}, S)$ for all $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in S$.
- $P_S(\mathbf{x} + \mathbf{y}) = P_S(\mathbf{x}) + \mathbf{y}$ for all $\mathbf{x} \in \mathbb{R}^n$ and $\mathbf{y} \in S$.
- $\text{dist}(a\mathbf{x}, S) = |a| \text{dist}(\mathbf{x}, S)$ for all $\mathbf{x} \in \mathbb{R}^n$ and real a .
- $P_S(a\mathbf{x}) = aP_S(\mathbf{x})$ for all $\mathbf{x} \in \mathbb{R}^n$ and real a .

7. Let \mathbf{M} be an $n \times n$ symmetric matrix with spectral decomposition $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{U}^*$, where \mathbf{U} is an orthogonal matrix and \mathbf{D} is a diagonal matrix with i th diagonal entry d_i . Prove that the Frobenius norm $\|\mathbf{M} - P_S(\mathbf{M})\|_F$ is minimized over the closed convex cone S of positive semidefinite matrices by taking $P_S(\mathbf{M}) = \mathbf{U}\mathbf{D}_+\mathbf{U}^*$, where \mathbf{D}_+ is diagonal with i th diagonal entry $\max\{d_i, 0\}$.

8. Let $C = \{(\mathbf{x}, t) \in \mathbb{R}^{n+1} : \|\mathbf{x}\| \leq t\}$ denote the ice cream cone in \mathbb{R}^{n+1} . Verify the projection formulas

$$P_C[(\mathbf{x}, t)] = \begin{cases} (\mathbf{x}, t) & \|\mathbf{x}\| \leq t \text{ and } t \geq 0 \\ (\mathbf{0}, 0) & \|\mathbf{x}\| \leq -t \text{ and } t \leq 0 \\ \left(\frac{\|\mathbf{x}\|+t}{2\|\mathbf{x}\|}\mathbf{x}, \frac{\|\mathbf{x}\|+t}{2}\right) & \text{otherwise.} \end{cases}$$

9. The convex regression example in Sect. 15.2 implicitly assumes that the regression function is defined on the integers $1, \dots, n$. Consider instead the problem of finding a convex function $f(\mathbf{x})$ on \mathbb{R}^m such that the sum of squares $\sum_{i=1}^n [y_i - f(\mathbf{x}_i)]^2$ is minimized. Demonstrate that this alternative problem can be rephrased as the quadratic programming problem of minimizing $\sum_{i=1}^n (y_i - z_i)^2$ subject to the restrictions

$$z_j \geq z_i + \mathbf{g}_i^*(\mathbf{x}_j - \mathbf{x}_i)$$

for all i and $j \neq i$. Here the unknown subgradients \mathbf{g}_i must be found along with the function values z_i at the specified points \mathbf{x}_i [19].

10. Verify the projection formula in Problem 11 of Chap. 5 by invoking the obtuse angle criterion.
11. Suppose C is a closed convex set wholly contained within an affine subspace $V = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{A}\mathbf{y} = \mathbf{b}\}$. For $\mathbf{x} \notin V$ demonstrate the projection identity $P_C(\mathbf{x}) = P_C \circ P_V(\mathbf{x})$ [196]. (Hint: Consider the equality

$$\begin{aligned} [\mathbf{x} - P_C(\mathbf{x})]^*[\mathbf{y} - P_C(\mathbf{x})] &= [\mathbf{x} - P_V(\mathbf{x})]^*[\mathbf{y} - P_C(\mathbf{x})] \\ &\quad + [P_V(\mathbf{x}) - P_C(\mathbf{x})]^*[\mathbf{y} - P_C(\mathbf{x})] \end{aligned}$$

with the obtuse angle criterion in mind.)

12. For positive numbers c_1, \dots, c_n and nonnegative numbers b_1, \dots, b_n satisfying $\sum_{i=1}^n c_i b_i \leq 1$, define the truncated simplex

$$S = \left\{ \mathbf{y} \in \mathbb{R}^n : \sum_{i=1}^n c_i y_i = 1, y_i \geq b_i, 1 \leq i \leq n \right\}.$$

If $\mathbf{x} \in \mathbb{R}^n$ has coordinate sum $\sum_{i=1}^n c_i x_i = 1$, then prove that the closest point \mathbf{y} in S to \mathbf{x} satisfies the Lagrange multiplier conditions

$$y_i - x_i + \lambda c_i - \mu_i = 0$$

for appropriate multipliers λ and $\mu_i \geq 0$. Further show that

$$\lambda = \frac{\mathbf{c}^* \boldsymbol{\mu}}{\|\mathbf{c}\|^2} \geq 0.$$

Why does it follow that $y_i = b_i$ whenever $x_i < b_i$? Prove that the Lagrange multiplier conditions continue to hold when $x_i < b_i$ if we replace x_i by b_i and μ_i by λc_i . Since the Lagrange multiplier conditions are sufficient as well as necessary in convex programming, this demonstrates that (a) we can replace each coordinate x_i by $\max\{x_i, b_i\}$ without changing the projection \mathbf{y} of \mathbf{x} onto S , and (b) \mathbf{y} can be viewed as a point in a similar simplex in a reduced number of dimensions when one or more $x_i \leq b_i$ [196].

13. Michelot's [196] algorithm for projecting a point \mathbf{x} onto the simplex S defined in Problem 12 cycles through the following steps:

- (a) Project onto the affine subspace $V_n = \{\mathbf{y} \in \mathbb{R}^n : \sum_i c_i y_i = 1\}$,
- (b) Replace each coordinate x_i by $\max\{x_i, b_i\}$,
- (c) Reduce the dimension n whenever some $x_i = b_i$.

In view of Problems 12 and 13, demonstrate that Michelot's algorithm converges to the correct solution in at most n steps. Explicitly solve the Lagrange multiplier problem corresponding to step (a). Program and test the algorithm.

14. Consider the problem of projecting a point \mathbf{x} onto the ℓ_1 ball

$$B = \{\mathbf{y} : \|\mathbf{y} - \mathbf{z}\|_1 \leq r\}.$$

Show that it suffices to project $\mathbf{x} - \mathbf{z}$ onto $\{\mathbf{w} : \|\mathbf{w}\|_1 \leq r\}$ and then translate the solution $\hat{\mathbf{w}}$ by \mathbf{z} . Hence, assume $\mathbf{z} = \mathbf{0}$ without loss of generality. Now argue that every entry of a solution $\hat{\mathbf{y}}$ should have the same sign as the corresponding entry of \mathbf{x} and that when $x_i = 0$, it does no harm to take $\hat{y}_i \geq 0$. Finally, sketch how projection onto an ℓ_1 ball can be achieved by projection onto the unit simplex.

15. The $\ell_{1,2}$ norm on \mathbb{R}^n is useful in group penalties [230]. Suppose the sets σ_g partition $\{1, \dots, n\}$ into groups with g as the group index. For $\mathbf{x} \in \mathbb{R}^n$ let \mathbf{x}_{σ_g} denote the vector formed by taking the components of \mathbf{x} derived from σ_g . The $\ell_{1,2}$ norm equals

$$\|\mathbf{x}\|_{1,2} = \sum_g \|\mathbf{x}_{\sigma_g}\|.$$

Check that $\|\mathbf{x}\|_{1,2}$ satisfies the properties of a norm. Now consider projecting a point \mathbf{x} onto the ball $B_r = \{\mathbf{y} : \|\mathbf{y}\|_{1,2} \leq r\}$. If we let $c_g = \|\mathbf{x}_{\sigma_g}\|$ and suppose that $\sum_g c_g \leq r$, then the solution is \mathbf{x} . Thus, assume the contrary. If any $c_g = \|\mathbf{x}_{\sigma_g}\| = 0$, argue that one should take $\mathbf{y}_{\sigma_g} = \mathbf{0}$. Assume therefore that $c_g > 0$ for every g . Collect the Euclidean distances c_g into a vector \mathbf{c} , and project \mathbf{c} onto the closest point \mathbf{d} in the ℓ_1 ball of radius r . This can be accomplished

by the algorithm described in Problem 14. Finally, prove that the projection \mathbf{y} of \mathbf{x} onto B_r satisfies $\mathbf{y}_{\sigma_g} = d_g c_g^{-1} \mathbf{x}_{\sigma_g}$. (Hints: Suppose that $r_g \geq 0$ for all g and $\sum_g r_g \leq r$. The problem of minimizing $\|\mathbf{y} - \mathbf{x}\|_{1,2}^2$ subject to $\|\mathbf{y}_{\sigma_g}\| \leq r_g$ for each g is separable and can be solved by projection onto the pertinent Euclidean balls. Substitution now leads to the secondary problem of minimizing

$$\sum_g (c_g - r_g)^2 1_{\{r_g \leq c_g\}} \quad (15.11)$$

subject to $r_g \geq 0$ for all g and $\sum_g r_g \leq r$. Suppose the optimal choice of the vector \mathbf{r} involves $r_g > c_g$ for some g . There must be a corresponding g' with $r_{g'} < c_{g'}$. One can decrease the criterion (15.11) by decreasing r_g and increasing $r_{g'}$. Hence, all $r_g \leq c_g$ at the optimal point, and one can dispense with the indicators $1_{\{r_g \geq c_g\}}$ and minimize the criterion (15.11) by ordinary projection.)

16. A polyhedral set is the nonempty intersection of a finite number of halfspaces. Program Dykstra's algorithm for projection onto the closest point of an arbitrary polyhedral set. Also program cyclic projection as suggested in Proposition 15.3.1, and compare it to Dykstra's algorithm on one or more test problems.
17. Demonstrate that the map $f(x) = x + e^{-x}$ is contractive on $[0, \infty)$ but lacks a fixed point. Is $f(x)$ strictly contractive?
18. Prove that the iteration scheme

$$x_{m+1} = \frac{1}{1 + x_m}$$

with domain $[0, \infty)$ has one fixed point y and that y is globally attractive. Is the function $f(x) = (1 + x)^{-1}$ contractive or strictly contractive?

19. Consider the map

$$T(\mathbf{x}) = \begin{pmatrix} \frac{1}{2(1+x_2)} \\ \frac{1}{2}e^{-x_1} \end{pmatrix}$$

from $\mathbb{R}_+^2 = \{\mathbf{x} \in \mathbb{R}^2 : \mathbf{x} \geq \mathbf{0}\}$ to itself. Calculate the differential

$$dT(\mathbf{x}) = \begin{pmatrix} 0 & -\frac{1}{2(1+x_2)^2} \\ -\frac{1}{2}e^{-x_1} & 0 \end{pmatrix},$$

and show that $\|dT(\mathbf{x})\| \leq \frac{1}{2}$ for all \mathbf{x} . Deduce the mean value inequality

$$\|T(\mathbf{y}) - T(\mathbf{z})\| \leq \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|$$

implying that $T(\mathbf{x})$ is a strict contraction with a unique fixed point.

20. Let $\mathbf{M} = \mathbf{U}\mathbf{D}\mathbf{U}^{-1}$ be an $n \times n$ diagonalizable matrix, where \mathbf{D} is diagonal and \mathbf{U} is invertible. Here the i th column \mathbf{u}_i of \mathbf{U} is an eigenvector of \mathbf{M} with eigenvalue d_i equal to the i th diagonal entry of \mathbf{D} . For a vector \mathbf{x} with expansion $\sum_{i=1}^n a_i \mathbf{u}_i$, define the vector norm $\|\mathbf{x}\|_{\dagger} = \sum_{i=1}^n |a_i|$. Verify that this defines a norm with induced matrix norm $\|\mathbf{M}\|_{\dagger} = \max_{1 \leq i \leq n} |d_i|$. Show that the affine map $\mathbf{x} \mapsto \mathbf{M}\mathbf{x} + \mathbf{v}$ is a contraction under $\|\mathbf{x}\|_{\dagger}$ whenever the spectral radius of \mathbf{M} is strictly less than 1. Note that the d_i and \mathbf{u}_i may be complex. What is the fixed point of the map?
21. Suppose the maps T_0, \dots, T_{r-1} are paracontractive with fixed point sets F_0, \dots, F_{r-1} . If $F = \bigcap_{i=0}^{r-1} F_i$ is nonempty, then show that the map $S = T_{r-1} \circ \dots \circ T_0$ is paracontractive with fixed point set $F = \bigcap_{i=0}^{r-1} F_i$.
22. Suppose the continuous map T from a closed convex set D to itself has a k -fold composition $S = T \circ \dots \circ T$ that is a strict contraction. Demonstrate that T and S share a unique fixed point, and that $\mathbf{x}_{m+1} = T(\mathbf{x}_m)$ converges to it.
23. Let $T(\mathbf{x})$ map the compact convex set C into itself. If there is a norm $\|\mathbf{x}\|_{\dagger}$ under which $\|T(\mathbf{y}) - T(\mathbf{x})\|_{\dagger} \leq \|\mathbf{y} - \mathbf{x}\|_{\dagger}$ for all \mathbf{y} and \mathbf{x} , then show that $T(\mathbf{x})$ has a fixed point. Use this result to prove that every finite state Markov chain possesses a stationary distribution. (Hints: Choose any $\mathbf{z} \in C$ and $\epsilon \in (0, 1)$ and define

$$T_{\epsilon}(\mathbf{x}) = (1 - \epsilon)T(\mathbf{x}) + \epsilon\mathbf{z}.$$

Argue that $T_{\epsilon}(\mathbf{x})$ is a strict contraction and send ϵ to 0.)

24. Under the hypotheses of Problem 23, demonstrate that the set of fixed points is nonempty, compact, and convex. (Hint: To prove convexity, suppose \mathbf{x} and \mathbf{y} are fixed points. For $\lambda \in [0, 1]$, argue that the point $\mathbf{z} = \lambda\mathbf{x} + (1 - \lambda)\mathbf{y}$ satisfies

$$\begin{aligned} \|\mathbf{y} - \mathbf{x}\|_{\dagger} &\leq \|T(\mathbf{y}) - T(\mathbf{z})\|_{\dagger} + \|T(\mathbf{z}) - T(\mathbf{x})\|_{\dagger} \\ &\leq \|\mathbf{y} - \mathbf{x}\|_{\dagger}. \end{aligned}$$

Deduce from this result that $T(\mathbf{z}) = \mathbf{z}$.)

25. Consider the problem of minimizing the convex function $f(\mathbf{x})$ subject to the affine equality constraints $g_i(\mathbf{x}) = 0$ for $1 \leq i \leq p$ and the convex inequality constraints $h_j(\mathbf{x}) \leq c_j$ for $1 \leq j \leq q$. Let $v(\mathbf{c})$ be the optimal value of $f(\mathbf{x})$ subject to the constraints. Show that the function $v(\mathbf{c})$ is convex in \mathbf{c} .
26. Describe the dual function for the power plan problem when one or more of the cost functions $f_i(x_i) = a_i x_i$ is linear instead of quadratic. Does this change affect the proposed solution by bisection or steepest ascent?

27. Calculate the dual function for the problem of minimizing $|x|$ subject to $x \leq -1$. Show that the optimal values of the primal and dual problems agree.
28. Verify that the two linear programs in Example 15.5.8 are dual programs.
29. Demonstrate that the dual function for Example 5.5.3 is

$$\mathcal{D}(\mu) = \begin{cases} 0 & \mu = 0 \\ 2 \sum_{i=1}^n \sqrt{\mu a_i c_i} - \mu b & \mu > 0. \end{cases}$$

Check that this problem is a convex program and that Slater's condition is satisfied.

30. Derive the dual function

$$\mathcal{D}(\lambda, \mu) = -\lambda - \mu^* \mathbf{e} - e^{-\lambda-1} \sum_{i=1}^n e^{-\mathbf{w}_i^* \mu}$$

for the problem of minimizing the negative entropy $\sum_{i=1}^n x_i \ln x_i$ subject to the constraints $\sum_{i=1}^n x_i = 1$, $\mathbf{W}\mathbf{x} \leq \mathbf{e}$, and all $x_i \geq 0$. Here the vector \mathbf{w}_i is the i th column of \mathbf{W} . (Hint: See Table 14.1.)

31. Consider the problem of minimizing the convex function

$$f(\mathbf{x}) = x_1 \ln x_1 - x_1 + x_2 \ln x_2 - x_2$$

subject to the constraints $x_1 + 2x_2 \leq 1$, $x_1 \geq 0$, and $x_2 \geq 0$. Show that the primal and dual optimal values coincide [17].

32. In the analytic centering program one minimizes the objective function $f(\mathbf{x}) = -\sum_{i=1}^n \ln x_i$ subject to the linear equality constraints $\mathbf{V}\mathbf{x} = \mathbf{d}$ and the positivity constraints $x_i > 0$ for all i . We can incorporate the positivity constraints into the objective function by defining $f(\mathbf{x}) = \infty$ if any $x_i \leq 0$. With this essential domain in mind, calculate the Fenchel conjugate

$$f^*(\mathbf{y}) = \begin{cases} -n - \sum_{i=1}^n \ln(-y_i) & \text{all } y_i < 0 \\ \infty & \text{any } y_i \geq 0. \end{cases}$$

Show that the dual problem consists in maximizing

$$-\lambda^* \mathbf{d} + n + \sum_{i=1}^n \ln(\mathbf{V}^* \lambda)_i$$

subject to the constraints $(\mathbf{V}^* \lambda)_i > 0$ for all i . The primal problem is easy to solve if we know the Lagrange multipliers. Indeed, straightforward differentiation shows that $x_i = 1/(\mathbf{V}^* \lambda)_i$. The moral here is

that if the number of rows of \mathbf{V} is small, then it is advantageous to solve the dual problem for λ and insert this value into the explicit solution for \mathbf{x} .

33. In the trust region problem of Sect. 11.7, one minimizes a quadratic $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$ subject to the constraint $\|\mathbf{x}\|^2 \leq r^2$. Why does the constraint satisfy Slater's condition? If \mathbf{A} is positive definite, then the problem is convex, and the primal and dual programs have the same optimal values. Calculate the dual function $\mathcal{D}(\mu)$ depending on the Lagrange multiplier μ . If \mathbf{A} is not positive definite, then show that

$$\frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + c = -\alpha r^2 + \frac{1}{2}\mathbf{x}^* (\mathbf{A} + \alpha \mathbf{I}) \mathbf{x} + \mathbf{b}^* \mathbf{x} + c$$

subject to the equality constraint $\|\mathbf{x}\|^2 = r^2$. How does this permit one to handle indefinite programs? The article [245] treats this problem in detail.

34. In the experimental design literature, the problems of D -optimal and A -optimal designs are well studied [225]. These involve minimizing the objective functions

$$f(\mathbf{x}) = \ln \det \left(\sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^* \right)^{-1}$$

$$g(\mathbf{x}) = \text{tr} \left(\sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^* \right)^{-1}$$

subject to the explicit constraints $\sum_{i=1}^p x_i = 1$ and $x_i \geq 0$ and the implicit constraint that the matrix $\sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^*$ is positive definite. The vectors $\mathbf{v}_1, \dots, \mathbf{v}_p$ in \mathbb{R}^n are given in advance. To formulate the dual problems, set $\mathbf{W} = \sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^*$ and require \mathbf{W} to be positive definite. Derive the dual problems with the constraint $\mathbf{W} = \sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^*$ included. Show that the dual problems involve maximizing $\ln \det \mathbf{X}$ and $\text{tr}(\mathbf{X}^{1/2})^2$, respectively, subject to the constraints that \mathbf{X} is positive definite and $\mathbf{v}_i^* \mathbf{X} \mathbf{v}_i \leq 1$ for all i . See the reference [262] for this formulation and techniques for maximizing the dual function. (Hints: The dual function for the objective function $f(\mathbf{x})$ is

$$\mathcal{D}(\mathbf{X}, \lambda) = \inf_{\mathbf{W}, \mathbf{x} \geq 0} \left\{ \ln \det \mathbf{W}^{-1} + \text{tr} \left[\mathbf{X} \left(\mathbf{W} - \sum_{i=1}^p x_i \mathbf{v}_i \mathbf{v}_i^* \right) \right] \right\}$$

$$+ \lambda \left(\sum_{i=1}^p x_i - 1 \right).$$

Eliminate λ by reparameterizing \mathbf{X} . For the objective function $g(\mathbf{x})$, take $p = 1$ in the next exercise.)

35. Let \mathbf{X} be a positive definite matrix and p a positive scalar. Demonstrate that the function $f(\mathbf{W}) = \text{tr}(\mathbf{W}^{-p}) + p \text{tr}(\mathbf{X}\mathbf{W})$ achieves its minimum of $(p+1) \text{tr}[\mathbf{X}^{p/(p+1)}]$ at the matrix $\mathbf{W} = \mathbf{X}^{-1/(p+1)}$. Here the argument \mathbf{W} is also assumed positive definite. (Hint: Reduce the problem to one of calculating a Fenchel conjugate via Proposition 14.6.1.)
36. We have repeatedly visited the problem of projecting an exterior point onto a closed convex set C . Consider a non-Euclidean norm $\|\mathbf{x}\|_{\dagger}$ with dual norm $\|\mathbf{y}\|_{\star}$. Let \mathbf{u} be a point exterior to C . Argue that the projection of \mathbf{u} onto C under this alternative norm can be found by minimizing the criterion $\|\mathbf{z}\|_{\dagger} + \delta_C(\mathbf{x})$ subject to $\mathbf{u} - \mathbf{x} = \mathbf{z}$. Demonstrate that the dual problem can be phrased as minimizing

$$\mathcal{D}(\boldsymbol{\lambda}) = \boldsymbol{\lambda}^* \mathbf{u} - \sup_{\mathbf{x} \in C} \boldsymbol{\lambda}^* \mathbf{x} \quad \text{for } \|\boldsymbol{\lambda}\|_{\star} \leq 1.$$

16

Convex Minimization Algorithms

16.1 Introduction

This chapter delves into three advanced algorithms for convex minimization. The projected gradient algorithm is useful in minimizing a strictly convex quadratic over a closed convex set. Although the algorithm extends to more general convex functions, the best theoretical results are available in this limited setting. We rely on the MM principle to motivate and extend the algorithm. The connections to Dykstra's algorithm and the contraction mapping principle add to the charm of the subject. On the minus side of the ledger, the projected gradient method can be very slow to converge. This defect is partially offset by ease of coding in many problems.

The second algorithm, path following in the exact penalty method, requires a fairly sophisticated understanding of convex calculus. As described in Chap. 13, classical penalty methods for solving constrained optimization problems exploit smooth penalties and send the tuning constant to infinity. If one substitutes absolute value and hinge penalties for square penalties, then there is no need to pass to the limit. Taking the penalty tuning constant sufficiently large generates a penalized problem with the same minimum as the constrained problem. In path following we track the minimum point of the penalized objective function as the tuning constant increases. Invocation of the implicit function theorem reduces path following to an exercise in numerically solving an ordinary differential equation [283, 284].

Our third algorithm, Bregman iteration [120, 208, 279], has found the majority of its applications in image processing. In ℓ_1 penalized image restoration, it gives sparser, better fitting signals. In total variation penalized image reconstruction, it gives higher contrast images with decent smoothing. The basis pursuit problem [38, 74] of minimizing $\|\mathbf{u}\|_1$ subject to $\mathbf{A}\mathbf{u} = \mathbf{f}$ readily succumbs to Bregman iteration. In many cases the basis pursuit solution is the sparsest consistent with the constraint. One can solve the basis pursuit problem by linear programming, but conventional solvers are not tailored to dense matrices \mathbf{A} and sparse solutions. Many applications require substitution of $\|\mathbf{D}\mathbf{u}\|_1$ for $\|\mathbf{u}\|_1$ for a smoothing matrix \mathbf{D} . This complication motivated the introduction of split Bregman iteration [106], which we briefly cover. Solution techniques continue to evolve rapidly in Bregman iteration. The whole field is driven by the realization that well-controlled sparsity gives better statistical inference and faster, more reliable algorithms than competing models and methods.

16.2 Projected Gradient Algorithm

For an $n \times n$ positive definite matrix \mathbf{A} and a closed convex set $S \subset \mathbb{R}^n$, consider the problem of minimizing the quadratic function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^*\mathbf{A}\mathbf{x} + \mathbf{b}^*\mathbf{x}$$

over S . We have studied this problem in depth in the special case $\mathbf{A} = \mathbf{I}$, where it reduces to projection onto S . Given that projection is relatively easy for a variety of convex sets, it is worth asking when the more general problem can be reduced to this special case. One way of answering the question is through the majorization

$$\begin{aligned} \mathbf{x}^*\mathbf{A}\mathbf{x} &= (\mathbf{x} - \mathbf{x}_m + \mathbf{x}_m)^*\mathbf{A}(\mathbf{x} - \mathbf{x}_m + \mathbf{x}_m) \\ &= (\mathbf{x} - \mathbf{x}_m)^*\mathbf{A}(\mathbf{x} - \mathbf{x}_m) + 2\mathbf{x}_m^*\mathbf{A}(\mathbf{x} - \mathbf{x}_m) + \mathbf{x}_m^*\mathbf{A}\mathbf{x}_m \\ &\leq \|\mathbf{A}\|_2\|\mathbf{x} - \mathbf{x}_m\|^2 + 2\mathbf{x}_m^*\mathbf{A}(\mathbf{x} - \mathbf{x}_m) + \mathbf{x}_m^*\mathbf{A}\mathbf{x}_m. \end{aligned}$$

This majorization leads to the function

$$g(\mathbf{x} \mid \mathbf{x}_m) = \frac{\|\mathbf{A}\|_2}{2}\|\mathbf{x} - \mathbf{x}_m\|^2 + \mathbf{x}_m^*\mathbf{A}(\mathbf{x} - \mathbf{x}_m) + \mathbf{b}^*(\mathbf{x} - \mathbf{x}_m) + c$$

majorizing $f(\mathbf{x})$, where c is an irrelevant constant. Completing the square allows one to rewrite the surrogate function as

$$g(\mathbf{x} \mid \mathbf{x}_m) = \frac{\|\mathbf{A}\|_2}{2} \left\| \mathbf{x} - \mathbf{x}_m + \frac{1}{\|\mathbf{A}\|_2}[\mathbf{A}\mathbf{x}_m + \mathbf{b}] \right\|^2 + d$$

for another irrelevant constant d . The majorization of $f(\mathbf{x})$ persists if we replace the coefficient $\|\mathbf{A}\|_2$ appearing in $g(\mathbf{x} \mid \mathbf{x}_m)$ by a larger constant.

The surrogate function $g(\mathbf{x} \mid \mathbf{x}_m)$ is essentially a Euclidean distance, and minimizing it over S is accomplished by projection onto S . If $P_S(\mathbf{y})$ is the projection operator, then the algorithm map boils down to

$$M(\mathbf{x}) = P_S \left[\mathbf{x} - \frac{1}{\|\mathbf{A}\|_2} (\mathbf{A}\mathbf{x} + \mathbf{b}) \right] = P_S \left[\mathbf{x} - \frac{1}{\|\mathbf{A}\|_2} \nabla f(\mathbf{x}) \right].$$

According to the MM principle, this projected gradient satisfies the descent property for the objective function $f(\mathbf{x})$. More generally, the algorithm map

$$M_\rho(\mathbf{x}) = P_S \left[\mathbf{x} - \frac{\rho}{\|\mathbf{A}\|_2} (\mathbf{A}\mathbf{x} + \mathbf{b}) \right] = P_S \left[\mathbf{x} - \frac{\rho}{\|\mathbf{A}\|_2} \nabla f(\mathbf{x}) \right].$$

also possesses the descent property for all $\rho \in (0, 1]$. Convergence of the projected gradient algorithm is guaranteed by the contraction mapping theorem stated in Proposition 15.3.2. Indeed, let $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$ denote the eigenvalues of \mathbf{A} . Since $\|P_S(\mathbf{u}) - P_S(\mathbf{v})\| \leq \|\mathbf{u} - \mathbf{v}\|$ for all points \mathbf{u} and \mathbf{v} and the matrix $\mathbf{I} - \alpha\mathbf{A}$ has i th eigenvalue $1 - \alpha\lambda_i$ for any constant α , we have

$$\begin{aligned} \|M_\rho(\mathbf{x}) - M_\rho(\mathbf{y})\| &\leq \left\| \mathbf{x} - \frac{\rho}{\|\mathbf{A}\|_2} (\mathbf{A}\mathbf{x} + \mathbf{b}) - \mathbf{y} + \frac{\rho}{\|\mathbf{A}\|_2} (\mathbf{A}\mathbf{y} + \mathbf{b}) \right\| \\ &= \left\| \left(\mathbf{I} - \frac{\rho}{\|\mathbf{A}\|_2} \mathbf{A} \right) (\mathbf{x} - \mathbf{y}) \right\| \\ &\leq \max_i \left| 1 - \frac{\rho\lambda_i}{\lambda_1} \right| \|\mathbf{x} - \mathbf{y}\|. \end{aligned}$$

Provided $\lambda_n > 0$ and $\rho \in (0, 2)$, it follows that the map $M_\rho(\mathbf{x})$ is a strict contraction on S . Except for a detail, this proves the second claim of the next proposition.

Proposition 16.2.1 *The projected gradient algorithm*

$$\mathbf{x}_{m+1} = P_S \left[\mathbf{x}_m - \frac{\rho}{\|\mathbf{A}\|_2} (\mathbf{A}\mathbf{x}_m + \mathbf{b}) \right] \quad (16.1)$$

for the convex quadratic $f(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A}\mathbf{x} + \mathbf{b}^* \mathbf{x}$ is a descent algorithm whenever $\rho \in (0, 1]$. The algorithm converges to the unique minimum of $f(\mathbf{x})$ on the convex set S whenever $f(\mathbf{x})$ is strictly convex and $\rho \in (0, 2)$.

Proof: Because the iteration map is a strict contraction, the iterates converge at a linear rate to its unique fixed point \mathbf{y} . It suffices to prove that \mathbf{y} furnishes the minimum. In view of the obtuse angle criterion, if \mathbf{z} is any point of S , we have

$$\left[\mathbf{y} - \frac{\rho}{\|\mathbf{A}\|_2} \nabla f(\mathbf{y}) - \mathbf{y} \right]^* [\mathbf{z} - \mathbf{y}] \leq 0.$$

However, this is equivalent to the condition $df(\mathbf{y})(\mathbf{z} - \mathbf{y}) \geq 0$, which is both necessary and sufficient for a minimum. ■

The special case of least squares estimation is important in Proposition 16.2.1. If

$$f(\mathbf{x}) = \frac{1}{2} \|\mathbf{y} - \mathbf{Q}\mathbf{x}\|^2 = \frac{1}{2} \mathbf{x}^* \mathbf{Q}^* \mathbf{Q} \mathbf{x} - \mathbf{y}^* \mathbf{Q} \mathbf{x} + \frac{1}{2} \|\mathbf{y}\|^2, \quad (16.2)$$

then in our previous notation $\mathbf{A} = \mathbf{Q}^* \mathbf{Q}$ and $\mathbf{b} = -\mathbf{Q}^* \mathbf{y}$. Furthermore, the matrix norm $\|\mathbf{A}\|_2 = \|\mathbf{Q}^* \mathbf{Q}\|_2 = \|\mathbf{Q}\|_2^2$. In practice finding the norm $\|\mathbf{A}\|_2$ is an issue. One can substitute the larger norm $\|\mathbf{A}\|_F$ for $\|\mathbf{A}\|_2$ as noted in Proposition 2.2.1. Alternatively, one can backtrack in the projected gradient algorithm (16.1). This involves making an initial guess a_0 of $\|\mathbf{A}\|_2$ and selecting a constant $c > 1$. If the point \mathbf{x}_{m+1} does not decrease $f(\mathbf{x})$, then replace a_0 by $a_1 = ca_0$ and recompute \mathbf{x}_{m+1} . If this new \mathbf{x}_{m+1} does not decrease $f(\mathbf{x})$, then replace a_1 by $a_2 = c^2 a_0$, and so forth. For k large enough, a_k must exceed $\|\mathbf{A}\|_2$. Problem 15 of Chap. 11 explores a simple algorithm of Hestenes and Karush [126] that efficiently produces the largest eigenvalue $\|\mathbf{A}\|_2$ of \mathbf{A} for n large.

Projected gradient algorithms are not limited to quadratic functions. Consider an arbitrary differentiable function $f(\mathbf{x})$ whose gradient satisfies the Lipschitz inequality

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq b \|\mathbf{y} - \mathbf{x}\|$$

for some b and all \mathbf{x} and \mathbf{y} . The projected gradient algorithm

$$\mathbf{x}_{m+1} = P_S \left[\mathbf{x}_m - \frac{\rho}{b} \nabla f(\mathbf{x}_m) \right]$$

with $\rho \in (0, 2)$ is designed to minimize $f(\mathbf{x})$ over a closed convex set S . Problems 1 and 2 sketch a few convergence results in this context [226]. See also Problem 31 of Chap. 4. For some functions the Lipschitz condition is only valid for a ball B centered at the current point \mathbf{x}_m . Then the algorithm that projects $\mathbf{x}_m - \frac{\rho}{b} \nabla f(\mathbf{x}_m)$ onto the intersection $B \cap S$ also retains the descent property. This amendment to the projected gradient method is inspired by the trust region strategy.

Example 16.2.1 *Projection onto the Image of a Convex Set*

Suppose projection onto the convex set S is easy. Given a compatible matrix \mathbf{Q} , the projected gradient algorithm allows us to project a point \mathbf{y} onto the image set $\mathbf{Q}S$. One merely minimizes the criterion (16.2) over S . For example, let S be the closed convex set $\{\mathbf{x} : \mathbf{x}_i \geq 0 \ \forall i > 1\}$, and let \mathbf{Q} be the lower-triangular matrix whose nonzero entries equal 1. The set $\mathbf{Q}S$ is the set $\{\mathbf{w} : w_1 \leq w_2 \leq \dots \leq w_n\}$ whose entries are nondecreasing.

The two transformations $\mathbf{w} = \mathbf{Q}\mathbf{x}$ and $\mathbf{v} = \mathbf{Q}^*\mathbf{u}$ can be implemented via the fast recurrences

$$\begin{aligned} w_1 &= x_1, & w_{k+1} &= w_k + x_{k+1} \\ v_n &= u_n, & v_k &= v_{k+1} + u_k. \end{aligned}$$

The first of these recurrences operates in a forward direction; the second operates in a backward direction. Projection onto $\mathbf{Q}S$ is quick because the recurrences and projection onto S are both quick. In practice, the pool adjacent violators algorithm of Example 15.2.3 is faster overall. ■

Example 16.2.2 Logistic Regression with a Group Lasso

Data from the National Opinion Research Center offers the opportunity to perform logistic regression with grouped categorical predictors [57]. Here we consider the dichotomy happy (very happy plus pretty happy) versus unhappy surveyed on $n = 1,566$ people. In addition to the primary response, each participant registered a level in five predictor categories: gender, marital status, education, financial status, and health. Within a category (or group) the first level is taken as baseline and omitted in analysis. In logistic regression we seek to maximize the loglikelihood

$$\begin{aligned} \ln L(\boldsymbol{\theta}) &= \sum_{i=1}^n [y_i \ln p_i + (1 - y_i) \ln(1 - p_i)] \\ y_i &\in \{0, 1\}, \quad p_i = \frac{e^{\mathbf{x}_i^* \boldsymbol{\theta}}}{1 + e^{\mathbf{x}_i^* \boldsymbol{\theta}}}. \end{aligned}$$

Here \mathbf{x}_i is the predictor vector for person i omitting the first level of each category; all entries of \mathbf{x}_i equal 0 or 1. The vector $\boldsymbol{\theta}$ encodes the corresponding regression coefficients. The ball $S = \{\boldsymbol{\theta} : \|\boldsymbol{\theta}\|_{1,2} \leq r\}$ is defined via the $\ell_{1,2}$ norm mentioned in Problem 15 of Chap. 15. This problem also outlines an effective algorithm for projection onto an $\ell_{1,2}$ ball.

Constrained maximization performs continuous model selection. The $\ell_{1,2}$ constraint groups the various parameters by category. As explained in Example 8.7, one can majorize $-\ln L(\boldsymbol{\theta})$ by the convex quadratic

$$f(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m) = -\ln L(\boldsymbol{\theta}_m) - d \ln L(\boldsymbol{\theta}_m)(\boldsymbol{\theta} - \boldsymbol{\theta}_m) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_m)^* \mathbf{A}(\boldsymbol{\theta} - \boldsymbol{\theta}_m),$$

where $\mathbf{A} = \frac{1}{4} \mathbf{X}^* \mathbf{X}$ and \mathbf{X} is the matrix whose i th row is \mathbf{x}_i^* . The projected gradient algorithm (16.1) with $\rho = 1$ therefore drives $f(\boldsymbol{\theta} \mid \boldsymbol{\theta}_m)$ downhill and $\ln L(\boldsymbol{\theta})$ uphill. Alternation of majorization and projection is effective in producing constrained maximum likelihood estimates in the limit.

Table 16.1 lists the estimated regression coefficients within each category as a function of the radius r . The most interesting findings in the table are the low impact of education on happiness and the parameter reversals

TABLE 16.1. Logistic regression with a Group-Lasso constraint

Radius r	0.00	0.20	0.40	0.60	0.80	1.00
Iterations	16	10	12	13	13	13
Female	–	–	–	–	–	–
Male	0.00	0.00	0.00	0.00	–0.00	–0.04
Married	–	–	–	–	–	–
Never married	0.00	0.00	0.00	0.00	–0.03	–0.06
Divorced	0.00	0.00	0.00	0.00	–0.05	–0.12
Widowed	0.00	0.00	0.00	0.00	–0.05	–0.12
Separated	0.00	0.00	0.00	0.00	–0.05	–0.12
Some high school	–	–	–	–	–	–
High school	0.00	0.00	0.00	0.00	0.00	0.00
Junior college	0.00	0.00	0.00	0.00	0.00	0.00
Bachelor	0.00	0.00	0.00	0.00	0.00	0.00
Graduate	0.00	0.00	0.00	0.00	0.00	0.00
Poor	–	–	–	–	–	–
Below average	0.00	–0.10	–0.16	–0.20	–0.21	–0.21
Average	0.00	0.08	0.13	0.18	0.20	0.21
Above average	0.00	0.07	0.13	0.20	0.23	0.25
Rich	0.00	0.01	0.01	0.02	0.03	0.03
Poor health	–	–	–	–	–	–
Fair health	0.00	–0.02	–0.06	–0.09	–0.09	–0.09
Good health	0.00	0.00	0.01	0.04	0.06	0.07
Excellent health	0.00	0.05	0.14	0.25	0.31	0.34

between poor and below-average financial status and between poor health and fair health. The number of iterations until convergence displayed in Table 16.1 suggest good numerical performance on this typical problem. ■

One can generalize the projected gradient algorithm in various ways. For instance, in the projected Newton method, one projects the partial Newton step $\mathbf{x}_m - \tau d^2 f(\mathbf{x}_m)^{-1} \nabla f(\mathbf{x}_m)$ onto the constraint set S [12]. Here the step length τ is usually taken to be 1, and the second differential $d^2 f(\mathbf{x}_m)$ is assumed positive definite. The projected Newton strategy tends to reduce the number of iterations while increasing the complexity per iteration. To minimize generic convex functions, one can substitute subgradients for gradients. Thus, one projects $\mathbf{x}_m - \tau \mathbf{g}_m$ onto S for $\mathbf{g}_m \in \partial f(\mathbf{x}_m)$ and some optimal choice of the constant τ . Unfortunately, there exist subgradients whose negatives are not descent directions. For instance, $-g$ is an ascent direction of $f(x) = |x|$ for any nontrivial $g \in \partial f(0) = [-1, 1]$. The next proposition partially salvages the situation.

Proposition 16.2.2 Suppose \mathbf{y} minimizes the convex function $f(\mathbf{x})$ over the closed convex set S . If $\mathbf{g}_m \in \partial f(\mathbf{x}_m)$,

$$\mathbf{x}_{m+1} = P_S[\mathbf{x}_m - \tau \mathbf{g}_m],$$

and \mathbf{x}_m is not optimal, then the choice

$$0 < \tau < \frac{2[f(\mathbf{x}_m) - f(\mathbf{y})]}{\|\mathbf{g}_m\|^2}$$

produces $\|\mathbf{x}_{m+1} - \mathbf{y}\| < \|\mathbf{x}_m - \mathbf{y}\|$.

Proof: Because projection is nonexpansive,

$$\begin{aligned} \|\mathbf{x}_{m+1} - \mathbf{y}\|^2 &= \|P_S(\mathbf{x}_m - \tau \mathbf{g}_m) - P_S(\mathbf{y})\|^2 \\ &\leq \|\mathbf{x}_m - \tau \mathbf{g}_m - \mathbf{y}\|^2 \\ &= \|\mathbf{x}_m - \mathbf{y}\|^2 - 2\tau \mathbf{g}_m^*(\mathbf{x}_m - \mathbf{y}) + \tau^2 \|\mathbf{g}_m\|^2. \end{aligned}$$

Hence, the inequality $f(\mathbf{y}) \geq f(\mathbf{x}_m) + \mathbf{g}_m^*(\mathbf{y} - \mathbf{x}_m)$ implies

$$\begin{aligned} \|\mathbf{x}_{m+1} - \mathbf{y}\|^2 &\leq \|\mathbf{x}_m - \mathbf{y}\|^2 + 2\tau[f(\mathbf{y}) - f(\mathbf{x}_m)] + \tau^2 \|\mathbf{g}_m\|^2 \\ &= \|\mathbf{x}_m - \mathbf{y}\|^2 + h(\tau) \end{aligned}$$

for the obvious quadratic $h(\tau)$, which satisfies $h(0) = 0$ and attains its minimum value

$$\min_{\tau} h(\tau) = -\frac{[f(\mathbf{x}_m) - f(\mathbf{y})]^2}{\|\mathbf{g}_m\|^2}$$

at the positive point

$$\hat{\tau} = \frac{f(\mathbf{x}_m) - f(\mathbf{y})}{\|\mathbf{g}_m\|^2}.$$

The claim now follows from the symmetry of $h(\tau)$ around the point $\hat{\tau}$. ■

In practice the value $f(\mathbf{y})$ is not known beforehand. This necessitates some strategy for choosing the step-length constants τ_m . For the sake of brevity, we refer the reader to the book [226] for further discussion.

16.3 Exact Penalties and Lagrangians

In nonlinear programming, exact penalty methods minimize the function

$$\mathcal{E}_{\rho}(\mathbf{y}) = f(\mathbf{y}) + \rho \sum_{i=1}^p |g_i(\mathbf{y})| + \rho \sum_{j=1}^q \max\{0, h_j(\mathbf{y})\},$$

where $f(\mathbf{y})$ is the objective function, $g_i(\mathbf{y})$ is an equality constraint, and $h_j(\mathbf{y})$ is an inequality constraint. It is interesting to compare this function to the Lagrangian function

$$\mathcal{L}(\mathbf{y}) = f(\mathbf{y}) + \sum_{i=1}^p \lambda_i g_i(\mathbf{y}) + \sum_{j=1}^q \mu_j h_j(\mathbf{y})$$

capturing the behavior of $f(\mathbf{y})$ near a constrained local minimum \mathbf{x} . Proposition 5.2.1 demonstrates that the Lagrangian satisfies the stationarity condition $\nabla \mathcal{L}(\mathbf{x}) = \mathbf{0}$; its inequality multipliers are nonnegative and obey the complementary slackness requirements $\mu_j h_j(\mathbf{x}) = 0$. In an exact penalty method we take

$$\rho > \max\{|\lambda_1|, \dots, |\lambda_p|, \mu_1, \dots, \mu_q\}. \quad (16.3)$$

This choice creates the favorable circumstances

$$\begin{aligned} \mathcal{L}(\mathbf{y}) &< \mathcal{E}_\rho(\mathbf{y}) \quad \text{for all infeasible } \mathbf{y} \\ \mathcal{L}(\mathbf{z}) &\leq f(\mathbf{z}) = \mathcal{E}_\rho(\mathbf{z}) \quad \text{for all feasible } \mathbf{z} \\ \mathcal{L}(\mathbf{x}) &= f(\mathbf{x}) = \mathcal{E}_\rho(\mathbf{x}) \quad \text{for } \mathbf{x} \text{ optimal} \end{aligned}$$

with profound consequences. As the next proposition proves, minimizing $\mathcal{E}_\rho(\mathbf{y})$ is effective in minimizing $f(\mathbf{y})$ subject to the constraints.

In most problems the Lagrange multipliers are unknown, so that a certain amount of trial and error is necessary to ensure that ρ is large enough. Alternatively, one can simply follow the exact solution path until it merges with the constrained solution. Before we study path following in detail, it is helpful to prove necessary and sufficient conditions for the two solutions to coincide. We start with the sufficient conditions expressed in Proposition 5.5.2.

Proposition 16.3.1 *Under the assumptions of Proposition 5.5.2, a constrained local minimum \mathbf{x} of $f(\mathbf{y})$ is an unconstrained local minimum of $\mathcal{E}_\rho(\mathbf{y})$ given inequality (16.3).*

Proof: Suppose the contrary is true, and let \mathbf{x}_m be a sequence of points that converge to \mathbf{x} and satisfy $\mathcal{E}_\rho(\mathbf{x}_m) < \mathcal{E}_\rho(\mathbf{x})$. Without loss of generality, assume that the unit vectors

$$\mathbf{v}_m = \frac{1}{\|\mathbf{x}_m - \mathbf{x}\|} (\mathbf{x}_m - \mathbf{x})$$

converge to a unit vector \mathbf{v} . Now consider the difference quotients

$$\frac{\mathcal{L}(\mathbf{x}_m) - \mathcal{L}(\mathbf{x})}{\|\mathbf{x}_m - \mathbf{x}\|} \leq \frac{\mathcal{E}_\rho(\mathbf{x}_m) - \mathcal{E}_\rho(\mathbf{x})}{\|\mathbf{x}_m - \mathbf{x}\|} < 0. \quad (16.4)$$

A brief calculation shows that

$$\begin{aligned} \lim_{m \rightarrow \infty} \frac{\mathcal{L}(\mathbf{x}_m) - \mathcal{L}(\mathbf{x})}{\|\mathbf{x}_m - \mathbf{x}\|} &= \lim_{m \rightarrow \infty} s_{\mathcal{L}}(\mathbf{x}_m, \mathbf{x}) \mathbf{v}_m \\ &= d\mathcal{L}(\mathbf{x}) \mathbf{v} \\ &= df(\mathbf{x}) \mathbf{v} + \sum_{i=1}^p \lambda_i dg_i(\mathbf{x}) \mathbf{v} + \sum_{j=1}^q \mu_j dh_j(\mathbf{x}) \mathbf{v}, \end{aligned}$$

where $s_{\mathcal{L}}(\mathbf{y}, \mathbf{x})$ is the slope function of $\mathcal{L}(\mathbf{y})$ around \mathbf{x} . (See Sect. 4.4 for a discussion of slope functions.) The stationarity condition $\nabla \mathcal{L}(\mathbf{x}) = \mathbf{0}$ implies $d\mathcal{L}(\mathbf{x}) \mathbf{v} = 0$.

The limit of the difference quotient for $\mathcal{E}_{\rho}(\mathbf{y})$ is more subtle to calculate. Under the subscript convention for slope functions, the equality $g_i(\mathbf{x}) = 0$ implies

$$\lim_{m \rightarrow \infty} \frac{|g_i(\mathbf{x}_m)| - |g_i(\mathbf{x})|}{\|\mathbf{x}_m - \mathbf{x}\|} = \lim_{m \rightarrow \infty} |s_{g_i}(\mathbf{x}_m, \mathbf{x}) \mathbf{v}_m| = |dg_i(\mathbf{x}) \mathbf{v}|.$$

The equality $h_j(\mathbf{x}) = 0$ for an active inequality constraint entails

$$\lim_{m \rightarrow \infty} \frac{\max\{0, h_j(\mathbf{x}_m)\} - \max\{0, h_j(\mathbf{x})\}}{\|\mathbf{x}_m - \mathbf{x}\|} = \max\{0, dh_j(\mathbf{x}) \mathbf{v}\}.$$

The inequalities $h_j(\mathbf{x}) < 0$ and $h_j(\mathbf{x}_m) < 0$ for an inactive inequality constraint likewise entail

$$\lim_{m \rightarrow \infty} \frac{\max\{0, h_j(\mathbf{x}_m)\} - \max\{0, h_j(\mathbf{x})\}}{\|\mathbf{x}_m - \mathbf{x}\|} = 0.$$

Hence, if the first r inequality constraints are active, then inequality (16.4) yields

$$\begin{aligned} 0 &\geq \lim_{m \rightarrow \infty} \frac{\mathcal{E}_{\rho}(\mathbf{x}_m) - \mathcal{E}_{\rho}(\mathbf{x})}{\|\mathbf{x}_m - \mathbf{x}\|} \\ &= df(\mathbf{x}) \mathbf{v} + \rho \sum_{i=1}^p |dg_i(\mathbf{x}) \mathbf{v}| + \rho \sum_{j=1}^r \max\{0, dh_j(\mathbf{x}) \mathbf{v}\}. \end{aligned}$$

Subtracting $d\mathcal{L}(\mathbf{x}) \mathbf{v} = 0$ from the last inequality gives

$$\begin{aligned} 0 &\geq \sum_{i=1}^p [\rho |dg_i(\mathbf{x}) \mathbf{v}| - \lambda_i dg_i(\mathbf{x}) \mathbf{v}] \\ &\quad + \sum_{j=1}^r [\rho \max\{0, dh_j(\mathbf{x}) \mathbf{v}\} - \mu_j dh_j(\mathbf{x}) \mathbf{v}] \\ &\geq \sum_{i=1}^p [\rho |dg_i(\mathbf{x}) \mathbf{v}| - \lambda_i dg_i(\mathbf{x}) \mathbf{v}] \\ &\quad + \sum_{j=1}^r (\rho - \mu_j) \max\{0, dh_j(\mathbf{x}) \mathbf{v}\}. \end{aligned}$$

Because all of the terms in the last two sums are nonnegative, they in fact vanish. But these are precisely the tangency conditions $dg_i(\mathbf{x})\mathbf{v} = 0$ and $dh_j(\mathbf{x})\mathbf{v} \leq 0$ for $h_j(\mathbf{x})$ active.

To finish the proof, we now pass to the limit in the second-order Taylor expansion

$$\frac{1}{2}\mathbf{v}_m^* s_{\mathcal{L}}^2(\mathbf{x}_m, \mathbf{x})\mathbf{v}_m = \frac{1}{\|\mathbf{x}_m - \mathbf{x}\|^2} [\mathcal{L}(\mathbf{x}_m) - \mathcal{L}(\mathbf{x})] < 0$$

and conclude that $\mathbf{v}^* d^2\mathcal{L}(\mathbf{x})\mathbf{v} \leq 0$, contrary to the assumption of Proposition 5.5.2 that no such unit tangent vector exists. Thus, the supposition that \mathbf{x} is not a local minimum of $\mathcal{E}_\rho(\mathbf{y})$ is untenable. ■

Further theoretical progress can be made by assuming that the equality constraints are affine and that the objective and inequality constraint functions are convex in addition to being differentiable. In these circumstances $\mathcal{E}_\rho(\mathbf{x})$ is convex owing to the closure properties of convex functions described in Proposition 6.3.3. Furthermore, at a feasible point \mathbf{x} with the first r inequality constraints active, the sum and chain rules yield

$$\partial\mathcal{E}_\rho(\mathbf{x}) = \nabla f(\mathbf{x}) + \rho \sum_{i=1}^p [-1, 1]\nabla g_i(\mathbf{x}) + \rho \sum_{j=1}^r [0, 1]\nabla h_j(\mathbf{x}).$$

Hence, if

$$\mathbf{0} = \nabla f(\mathbf{x}) + \sum_{i=1}^p \lambda_i \nabla g_i(\mathbf{x}) + \sum_{j=1}^r \mu_j \nabla h_j(\mathbf{x}), \quad (16.5)$$

then certainly $\mathbf{0} \in \partial\mathcal{E}_\rho(\mathbf{x})$. Thus, a constrained minimum point of $f(\mathbf{y})$ satisfying the multiplier rule (16.5) corresponds to an unconstrained minimum point of $\mathcal{E}_\rho(\mathbf{y})$. This is just the content of Proposition 16.3.1 specialized to convex programming.

Conversely, assume that \mathbf{z} is an unconstrained minimum point of $\mathcal{E}_\rho(\mathbf{y})$ and that \mathbf{x} is a constrained minimum point of $f(\mathbf{y})$ satisfying the multiplier rule (16.5). If \mathbf{z} is feasible, then \mathbf{z} also furnishes a constrained minimum of $f(\mathbf{y})$ because $f(\mathbf{y})$ and $\mathcal{E}_\rho(\mathbf{y})$ coincide on the feasible region. If \mathbf{z} is infeasible, then

$$\mathcal{E}_\rho(\mathbf{x}) = f(\mathbf{x}) = \mathcal{L}(\mathbf{x}) \leq \mathcal{L}(\mathbf{z}) < \mathcal{E}_\rho(\mathbf{z}).$$

Here the inequality $\mathcal{L}(\mathbf{x}) \leq \mathcal{L}(\mathbf{z})$ reflects the fact that the convex function $\mathcal{L}(\mathbf{y})$ attains its minimum at the stationary point \mathbf{x} . The contradiction $\mathcal{E}_\rho(\mathbf{x}) < \mathcal{E}_\rho(\mathbf{z})$ now shows that \mathbf{z} is feasible and consequently furnishes a constrained minimum of $f(\mathbf{y})$. For the sake of completeness, we restate this result as a formal proposition.

Proposition 16.3.2 *Suppose in a convex program with differentiable objective and constraint functions that there exists a constrained minimum \mathbf{x} of $f(\mathbf{y})$ satisfying the multiplier rule. Under condition (16.3), a point \mathbf{z} furnishes a constrained minimum of $f(\mathbf{y})$ if and only if it furnishes an unconstrained minimum of $\mathcal{E}_\rho(\mathbf{y})$.*

Proof: See the forgoing discussion. ■

In path following, one tracks the postulated minimum point $\mathbf{x}(\rho)$ of $\mathcal{E}_\rho(\mathbf{y})$ as a function of ρ until ρ exceeds the Lagrange multiplier threshold. Thus, specifying a stationarity condition for $\mathcal{E}_\rho(\mathbf{y})$ is crucial. Unfortunately, our previous derivations of stationarity conditions assumed either differentiability or convexity. In general nonlinear programs, $\mathcal{E}_\rho(\mathbf{y})$ is neither. Here we tackle the stationarity condition via forward directional derivatives. Once again assume that \mathbf{x} is a local minimum of $\mathcal{E}_\rho(\mathbf{y})$. If the objective and constraint functions are differentiable at \mathbf{x} , then they possess directional derivatives at \mathbf{x} , and

$$d_{\mathbf{v}}\mathcal{E}_\rho(\mathbf{x}) = df(\mathbf{x})\mathbf{v} + \rho \sum_{i=1}^p d_{\mathbf{v}}|g_i(\mathbf{x})| + \rho \sum_{j=1}^q d_{\mathbf{v}} \max\{0, h_j(\mathbf{x})\}.$$

In proving Proposition 16.3.1, we calculated the obscure pieces making up $d_{\mathbf{v}}\mathcal{E}_\rho(\mathbf{y})$. Based on the notation

$$\begin{aligned} \mathcal{N}_E &= \{i : g_i(\mathbf{x}) < 0\} & \mathcal{N}_I &= \{j : h_j(\mathbf{x}) < 0\} \\ \mathcal{Z}_E &= \{i : g_i(\mathbf{x}) = 0\} & \mathcal{Z}_I &= \{j : h_j(\mathbf{x}) = 0\} \\ \mathcal{P}_E &= \{i : g_i(\mathbf{x}) > 0\} & \mathcal{P}_I &= \{j : h_j(\mathbf{x}) > 0\} \end{aligned} \tag{16.6}$$

we have

$$\begin{aligned} d_{\mathbf{v}}\mathcal{E}_\rho(\mathbf{x}) &= df(\mathbf{x})\mathbf{v} - \rho \sum_{i \in \mathcal{N}_E} dg_i(\mathbf{x})\mathbf{v} + \rho \sum_{i \in \mathcal{P}_E} dg_i(\mathbf{x})\mathbf{v} + \rho \sum_{j \in \mathcal{P}_I} dh_j(\mathbf{x})\mathbf{v} \\ &\quad + \rho \sum_{i \in \mathcal{Z}_E} |dg_i(\mathbf{x})\mathbf{v}| + \rho \sum_{j \in \mathcal{Z}_I} \max\{0, dh_j(\mathbf{x})\mathbf{v}\} \\ &= \mathbf{w}^*\mathbf{v} + \rho \sum_{i \in \mathcal{Z}_E} |dg_i(\mathbf{x})\mathbf{v}| + \rho \sum_{j \in \mathcal{Z}_I} \max\{0, dh_j(\mathbf{x})\mathbf{v}\} \end{aligned}$$

for the obvious choice of \mathbf{w} . At a local minimum \mathbf{x} of $\mathcal{E}_\rho(\mathbf{y})$, all directional derivatives satisfy $d_{\mathbf{v}}\mathcal{E}_\rho(\mathbf{x}) \geq 0$.

We now focus on the function $\mathcal{K}(\mathbf{v}) = d_{\mathbf{v}}\mathcal{E}_\rho(\mathbf{x})$ and derive an appropriate stationarity condition. Since the composition of a convex function with a linear function is convex, $\mathcal{K}(\mathbf{v})$ is convex even when $\mathcal{E}_\rho(\mathbf{x})$ is not. Hence, in dealing with $\mathcal{K}(\mathbf{v})$, we can invoke the rules of the convex calculus developed in Sects. 14.4 and 14.5. Because $\mathcal{K}(\mathbf{v})$ achieves its minimum value of 0 at the origin $\mathbf{0}$, Proposition 14.4.3 implies the containment $\mathbf{0} \in \partial\mathcal{K}(\mathbf{0})$. Applying

the rules of convex differentiation to $\mathcal{K}(\mathbf{v})$ gives the subdifferential

$$\partial\mathcal{K}(\mathbf{0}) = \mathbf{w} + \rho \sum_{i \in \mathcal{Z}_E} [-1, 1] \nabla g_i(\mathbf{x}) + \rho \sum_{j \in \mathcal{Z}_I} [0, 1] \nabla h_j(\mathbf{x}).$$

The stationarity condition

$$\begin{aligned} \mathbf{0} \in & \nabla f(\mathbf{x}) - \rho \sum_{i \in \mathcal{N}_E} \nabla g_i(\mathbf{x}) + \rho \sum_{i \in \mathcal{P}_E} \nabla g_i(\mathbf{x}) + \rho \sum_{j \in \mathcal{P}_I} \nabla h_j(\mathbf{x}) \\ & + \rho \sum_{i \in \mathcal{Z}_E} [-1, 1] \nabla g_i(\mathbf{x}) + \rho \sum_{j \in \mathcal{Z}_I} [0, 1] \nabla h_j(\mathbf{x}) \end{aligned} \quad (16.7)$$

for $\mathcal{E}_\rho(\mathbf{x})$ is a consequence of these considerations.

16.4 Mechanics of Path Following

Throughout this section we restrict our attention to convex programs. As a prelude to our derivation of the path following algorithm, we record several properties of $\mathcal{E}_\rho(\mathbf{x})$ that mitigate the failure of differentiability.

Proposition 16.4.1 *The surrogate function $\mathcal{E}_\rho(\mathbf{x})$ is increasing in ρ . Furthermore, $\mathcal{E}_\rho(\mathbf{x})$ is strictly convex for one $\rho > 0$ if and only if it is strictly convex for all $\rho > 0$. Likewise, when $f(\mathbf{x})$ is bounded below, $\mathcal{E}_\rho(\mathbf{x})$ is coercive for one $\rho > 0$ if and only if it is coercive for all $\rho > 0$. Finally, if $f(\mathbf{x})$ is strictly convex (or coercive), then all $\mathcal{E}_\rho(\mathbf{x})$ are strictly convex (or coercive).*

Proof: The first assertion is obvious. For the second assertion, consider more generally a finite family $u_1(\mathbf{x}), \dots, u_q(\mathbf{x})$ of convex functions, and suppose a linear combination $\sum_{k=1}^q c_k u_k(\mathbf{x})$ with positive coefficients is strictly convex. It suffices to prove that any other linear combination $\sum_{k=1}^q b_k u_k(\mathbf{x})$ with positive coefficients is strictly convex. For any two points $\mathbf{x} \neq \mathbf{y}$ and any scalar $\alpha \in (0, 1)$, we have

$$u_k[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] \leq \alpha u_k(\mathbf{x}) + (1 - \alpha)u_k(\mathbf{y}). \quad (16.8)$$

Since $\sum_{k=1}^q c_k u_k(\mathbf{x})$ is strictly convex, strict inequality must hold for at least one k . Hence, multiplying inequality (16.8) by b_k and adding gives

$$\sum_{k=1}^q b_k u_k[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] < \alpha \sum_{k=1}^q b_k u_k(\mathbf{x}) + (1 - \alpha) \sum_{k=1}^q b_k u_k(\mathbf{y}).$$

The third assertion follows from the criterion given in Proposition 12.3.1. Indeed, suppose $\mathcal{E}_\rho(\mathbf{x})$ is coercive, but $\mathcal{E}_{\rho^*}(\mathbf{x})$ is not coercive. Then there exists a point \mathbf{x} , a direction \mathbf{v} , and a sequence of scalars t_n tending to

∞ such that $\mathcal{E}_{\rho^*}(\mathbf{x} + t_n \mathbf{v})$ is bounded above. This requires the sequence $f(\mathbf{x} + t_n \mathbf{v})$ and each of the sequences $|g_i(\mathbf{x} + t_n \mathbf{v})|$ and $\max\{0, h_j(\mathbf{x} + t_n \mathbf{v})\}$ to remain bounded above. But in this circumstance the sequence $\mathcal{E}_{\rho}(\mathbf{x} + t_n \mathbf{v})$ also remains bounded above. The final assertion is also obvious. ■

To speak coherently of solution paths, one must validate the existence, uniqueness, and continuity of the solution $\mathbf{x}(\rho)$ to the stationarity condition (16.7). Uniqueness follows from assuming that $f(\mathbf{x})$ is strictly convex or more generally by assuming that $\mathcal{E}_{\rho}(\mathbf{x})$ is strictly convex for a single positive ρ . Existence and continuity are more subtle. Let us restate the stationarity condition as

$$\mathbf{0} = \nabla f(\mathbf{x}) + \rho \sum_{i=1}^r s_i \nabla g_i(\mathbf{x}) + \rho \sum_{j=1}^s t_j \nabla h_j(\mathbf{x}) \quad (16.9)$$

for coefficient sets $\{s_i\}_{i=1}^r$ and $\{t_j\}_{j=1}^s$ that satisfy

$$s_i \in \begin{cases} \{-1\} & g_i(\mathbf{x}) < 0 \\ [-1, 1] & g_i(\mathbf{x}) = 0 \\ \{1\} & g_i(\mathbf{x}) > 0 \end{cases} \quad \text{and} \quad t_j \in \begin{cases} \{0\} & h_j(\mathbf{x}) < 0 \\ [0, 1] & h_j(\mathbf{x}) = 0 \\ \{1\} & h_j(\mathbf{x}) > 0. \end{cases} \quad (16.10)$$

This notation puts us into position to state and prove some basic facts.

Proposition 16.4.2 *If $\mathcal{E}_{\rho}(\mathbf{y})$ is strictly convex and coercive, then the solution path $\mathbf{x}(\rho)$ of equation (16.7) exists and is continuous in ρ . If the gradient vectors $\{\nabla g_i(\mathbf{x}) : g_i(\mathbf{x}) = 0\} \cup \{\nabla h_j(\mathbf{x}) : h_j(\mathbf{x}) = 0\}$ of the active constraints are linearly independent at $\mathbf{x}(\rho)$ for $\rho > 0$, then in addition the coefficients $s_i(\rho)$ and $t_j(\rho)$ are unique and continuous near ρ .*

Proof: In accord with Proposition 16.4.1, we assume that either $f(\mathbf{x})$ is strictly convex and coercive or restrict our attention to the open interval $(0, \infty)$. Consider a subinterval $[a, b]$ containing ρ and fix a point \mathbf{x} in the common domain of the functions $\mathcal{E}_{\rho}(\mathbf{y})$. The coercivity of $\mathcal{E}_a(\mathbf{y})$ and the inequalities

$$\mathcal{E}_a[\mathbf{x}(\rho)] \leq \mathcal{E}_{\rho}[\mathbf{x}(\rho)] \leq \mathcal{E}_{\rho}(\mathbf{x}) \leq \mathcal{E}_b(\mathbf{x})$$

demonstrate that the solution vector $\mathbf{x}(\rho)$ is bounded over $[a, b]$. To prove continuity, suppose that it fails for a given $\rho \in [a, b]$. Then there exists an $\epsilon > 0$ and a sequence ρ_n tending to ρ such $\|\mathbf{x}(\rho_n) - \mathbf{x}(\rho)\| \geq \epsilon$ for all n . Since $\mathbf{x}(\rho_n)$ is bounded, we can pass to a subsequence if necessary and assume that $\mathbf{x}(\rho_n)$ converges to some point \mathbf{y} . Taking limits in the inequality $\mathcal{E}_{\rho_n}[\mathbf{x}(\rho_n)] \leq \mathcal{E}_{\rho_n}(\mathbf{x})$ shows that $\mathcal{E}_{\rho}(\mathbf{y}) \leq \mathcal{E}_{\rho}(\mathbf{x})$ for all \mathbf{x} . Because $\mathbf{x}(\rho)$ is unique, we reach the contradictory conclusions $\|\mathbf{y} - \mathbf{x}(\rho)\| \geq \epsilon$ and $\mathbf{y} = \mathbf{x}(\rho)$.

Verification of the second claim is deferred to permit further discussion of path following. The claim says that an active constraint ($g_i(\mathbf{x}) = 0$ or $h_j(\mathbf{x}) = 0$) remains active until its coefficient hits an endpoint of its

subdifferential. Because the solution path is, in fact, piecewise smooth, one can follow the coefficient path by numerically solving an ordinary differential equation (ODE). ■

Along the solution path we keep track of the index sets defined in equation (16.6) and determined by the signs of the constraint functions. For the sake of simplicity, assume that at the beginning of the current segment s_i does not equal -1 or 1 when $i \in \mathcal{Z}_E$ and t_j does not equal 0 or 1 when $j \in \mathcal{Z}_I$. In other words, the coefficients of the active constraints occur on the interiors, either $(-1, 1)$ or $(0, 1)$, of their subdifferentials. Let us show in this circumstance that the solution path can be extended in a smooth fashion. Our plan of attack is to reparameterize by the Lagrange multipliers of the active constraints. Thus, set $\lambda_i = \rho s_i$ for $i \in \mathcal{Z}_E$ and $\omega_j = \rho t_j$ for $j \in \mathcal{Z}_I$. These multipliers satisfy $-\rho < \lambda_i < \rho$ and $0 < \omega_j < \rho$. The stationarity condition now reads

$$\begin{aligned} \mathbf{0} = & \nabla f(\mathbf{x}) - \rho \sum_{i \in \mathcal{N}_E} \nabla g_i(\mathbf{x}) + \rho \sum_{i \in \mathcal{P}_E} \nabla g_i(\mathbf{x}) + \rho \sum_{j \in \mathcal{P}_I} \nabla h_j(\mathbf{x}) \\ & + \sum_{i \in \mathcal{Z}_E} \lambda_i \nabla g_i(\mathbf{x}) + \sum_{j \in \mathcal{Z}_I} \omega_j \nabla h_j(\mathbf{x}). \end{aligned}$$

For convenience now define

$$\begin{aligned} \mathbf{U}_{\mathcal{Z}}(\mathbf{x}) &= \begin{bmatrix} dg_{\mathcal{Z}_E}(\mathbf{x}) \\ dh_{\mathcal{Z}_I}(\mathbf{x}) \end{bmatrix} \\ \mathbf{u}_{\bar{\mathcal{Z}}}(\mathbf{x}) &= - \sum_{i \in \mathcal{N}_E} \nabla g_i(\mathbf{x}) + \sum_{i \in \mathcal{P}_E} \nabla g_i(\mathbf{x}) + \sum_{j \in \mathcal{P}_I} \nabla h_j(\mathbf{x}). \end{aligned}$$

In this notation the stationarity equation can be recast as

$$\mathbf{0} = \nabla f(\mathbf{x}) + \rho \mathbf{u}_{\bar{\mathcal{Z}}}(\mathbf{x}) + \mathbf{U}_{\mathcal{Z}}^*(\mathbf{x}) \begin{pmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\omega} \end{pmatrix}.$$

Under the assumption that the matrix $\mathbf{U}_{\mathcal{Z}}(\mathbf{x})$ has full row rank, one can solve for the Lagrange multipliers in the form

$$\begin{pmatrix} \boldsymbol{\lambda} \\ \boldsymbol{\omega} \end{pmatrix} = -[\mathbf{U}_{\mathcal{Z}}(\mathbf{x})\mathbf{U}_{\mathcal{Z}}^*(\mathbf{x})]^{-1}\mathbf{U}_{\mathcal{Z}}(\mathbf{x})[\nabla f(\mathbf{x}) + \rho\mathbf{u}_{\bar{\mathcal{Z}}}(\mathbf{x})]. \quad (16.11)$$

Hence, the multipliers are unique. Continuity of the multipliers is a consequence of the continuity of the solution path $\mathbf{x}(\rho)$ and the continuity of all functions in sight on the right-hand side of equation (16.11). This observation completes the proof of Proposition 16.4.2.

In addition to the stationarity condition, one must enforce the constraint equations $0 = g_i(\mathbf{x})$ for $i \in \mathcal{Z}_E$ and $0 = h_j(\mathbf{x})$ for $j \in \mathcal{Z}_I$. Collectively the stationarity and constraint equations can be written as a vector equation $\mathbf{0} = \mathbf{k}(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}, \rho)$ with the active constraints appended below the stationarity condition. To solve for \mathbf{x} , $\boldsymbol{\lambda}$ and $\boldsymbol{\omega}$ in terms of ρ , we apply the implicit

function theorem. This requires calculating the differential of $k(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}, \rho)$ with respect to the underlying dependent variables \mathbf{x} , $\boldsymbol{\lambda}$, and $\boldsymbol{\omega}$ and the independent variable ρ . Because the equality constraints are affine, a brief calculation gives

$$\begin{aligned} \partial_{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}} k &= \begin{bmatrix} d^2 f(\mathbf{x}) + \rho \sum_{j \in \mathcal{P}_1} d^2 h_j(\mathbf{x}) + \sum_{j \in \mathcal{Z}_1} \omega_j d^2 h_j(\mathbf{x}) & \mathbf{U}_{\mathcal{Z}}^*(\mathbf{x}) \\ \mathbf{U}_{\mathcal{Z}}(\mathbf{x}) & \mathbf{0} \end{bmatrix} \\ \partial_{\rho} k &= \begin{bmatrix} \mathbf{u}_{\bar{\mathcal{Z}}}(\mathbf{x}) \\ \mathbf{0} \end{bmatrix}. \end{aligned}$$

In view of Proposition 5.2.2, the matrix $\partial_{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}} k(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}, \rho)$ is nonsingular when its upper-left block is positive definite and its lower-left block has full row rank. Given that it is nonsingular, the implicit function theorem applies, and we can in principle solve for \mathbf{x} , $\boldsymbol{\lambda}$ and $\boldsymbol{\omega}$ in terms of ρ . More importantly, the implicit function theorem supplies the derivative

$$\frac{d}{d\rho} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \\ \boldsymbol{\omega} \end{bmatrix} = -(\partial_{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}} k)^{-1} \partial_{\rho} k, \tag{16.12}$$

which is the key to path following. We summarize our findings in the next proposition.

Proposition 16.4.3 *Suppose the surrogate function $\mathcal{E}_{\rho}(\mathbf{y})$ is strictly convex and coercive. If at the point ρ_0 the matrix $\partial_{\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}} k(\mathbf{x}, \boldsymbol{\lambda}, \boldsymbol{\omega}, \rho)$ is nonsingular and the coefficient of each active constraint occurs on the interior of its subdifferential, then the solution path $\mathbf{x}(\rho)$ and Lagrange multipliers $\boldsymbol{\lambda}(\rho)$ and $\boldsymbol{\omega}(\rho)$ satisfy the differential equation (16.12) in the vicinity of ρ_0 .*

If one views ρ as time, then one can trace the solution path along the current time segment until either an inactive constraint becomes active or the coefficient of an active constraint hits the boundary of its subdifferential. The earliest hitting time or escape time over all constraints determines the duration of the current segment. When the hitting time for an inactive constraint occurs first, we move the constraint to the appropriate active set \mathcal{Z}_E or \mathcal{Z}_I and keep the other constraints in place. Similarly, when the escape time for an active constraint occurs first, we move the constraint to the appropriate inactive set and keep the other constraints in place. In the second scenario, if s_i hits the value -1 , then we move i to \mathcal{N}_E ; if s_i hits the value 1 , then we move i to \mathcal{P}_E . Similar comments apply when a coefficient t_j hits 0 or 1 . Once this move is executed, we commence path following along the new segment. Path following continues until for sufficiently large ρ , the sets \mathcal{N}_E , \mathcal{P}_E , and \mathcal{P}_1 are exhausted, $\mathbf{u}_{\bar{\mathcal{Z}}} = \mathbf{0}$, and the solution vector $\mathbf{x}(\rho)$ stabilizes.

Path following simplifies considerably in convex quadratic programming with objective function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x}$ and equality constraints

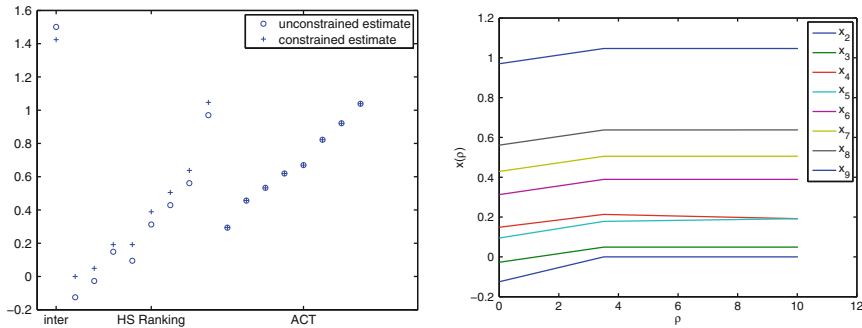


FIGURE 16.1. *Left:* Unconstrained and constrained estimates for the Iowa GPA data. *Right:* Solution paths for the high school rank regression coefficients

$\mathbf{V}\mathbf{x} = \mathbf{d}$ and inequality constraints $\mathbf{W}\mathbf{x} \leq \mathbf{e}$, where \mathbf{A} is positive semi-definite. The exact penalized objective function becomes

$$\mathcal{E}_\rho(\mathbf{x}) = \frac{1}{2}\mathbf{x}^* \mathbf{A} \mathbf{x} + \mathbf{b}^* \mathbf{x} + \rho \sum_{i=1}^s |\mathbf{v}_i^* \mathbf{x} - d_i| + \rho \sum_{j=1}^* (\mathbf{w}_j^* \mathbf{x} - e_j)_+.$$

Since both the equality and inequality constraints are affine, their second derivatives vanish. Both \mathbf{U}_Z and $\mathbf{u}_{\bar{z}}$ are constant on the current path segment, and the path $\mathbf{x}(\rho)$ satisfies

$$\frac{d}{d\rho} \begin{bmatrix} \mathbf{x} \\ \boldsymbol{\lambda} \\ \boldsymbol{\omega} \end{bmatrix} = - \begin{pmatrix} \mathbf{A} & \mathbf{U}_Z^* \\ \mathbf{U}_Z & \mathbf{0} \end{pmatrix}^{-1} \begin{pmatrix} \mathbf{u}_{\bar{z}} \\ \mathbf{0} \end{pmatrix}. \quad (16.13)$$

Because the solution path $\mathbf{x}(\rho)$ is piecewise linear, it is possible to anticipate the next hitting or exit time and take a large jump. The matrix inverse appearing in equation (16.13) can be efficiently updated by the sweep operator of computational statistics [283].

Example 16.4.1 Partial Isotone Regression

Order-constrained regression is now widely accepted as an important modeling tool in statistics [219, 236]. If \mathbf{x} is the parameter vector, monotone regression includes isotone constraints $x_1 \leq x_2 \leq \dots \leq x_m$ and antitone constraints $x_1 \geq x_2 \geq \dots \geq x_m$. In partially ordered regression, subsets of the parameters are subject to isotone or antitone constraints. As an example of partial isotone regression, consider the data from Table 1.3.1 of the reference [219] on the first-year grade point averages (GPA) of 2397 University of Iowa freshmen. These data can be downloaded as part of the R package `ic.infer`. The ordinal predictors high school rank (as a percentile) and ACT (a standard aptitude test) score are discretized into nine ordered categories each. It is rational to assume that college performance is isotone separately within each predictor set. Figure 16.1 shows

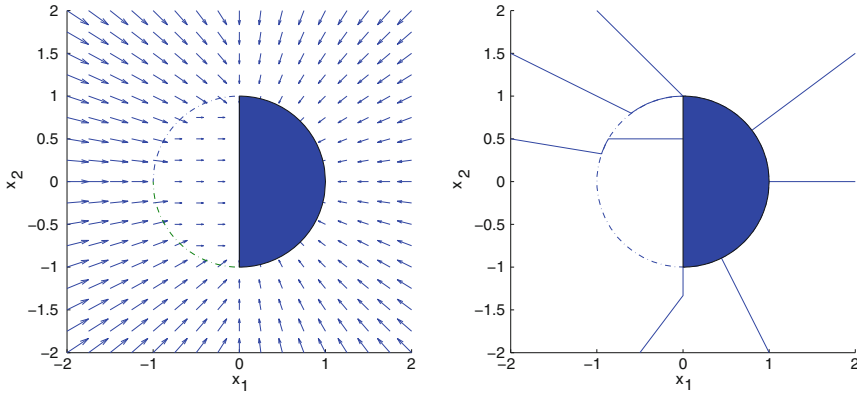


FIGURE 16.2. Projection to the positive half-disc. *Left:* Derivatives at $\rho = 0$ for projection onto the half-disc. *Right:* Projection trajectories from various initial points

the unconstrained and constrained solutions for the intercept and the two predictor sets and the solution path of the regression coefficients for the high school rank predictor. In this quadratic programming problem, the solution path is piecewise linear. In contrast the next example involves nonlinear path segments. ■

Example 16.4.2 *Projection onto the Half-Disc*

Dijkstra’s algorithm as explained in Sect. 15.2 projects an exterior point onto the intersection of a finite number of closed convex sets. The projection problem also yields to path following. Consider our previous toy example of projecting a point $\mathbf{b} \in \mathbb{R}^2$ onto the intersection of the closed unit ball and the closed half space $x_1 \geq 0$. This is equivalent to solving

$$\begin{aligned} \text{minimize} \quad & f(\mathbf{x}) = \frac{1}{2} \|\mathbf{x} - \mathbf{b}\|^2 \\ \text{subject to} \quad & h_1(\mathbf{x}) = \frac{1}{2} \|\mathbf{x}\|^2 - \frac{1}{2} \leq 0, \quad h_2(\mathbf{x}) = -x_1 \leq 0 \end{aligned}$$

with gradients and second differentials

$$\begin{aligned} \nabla f(\mathbf{x}) &= \mathbf{x} - \mathbf{b}, & \nabla h_1(\mathbf{x}) &= \mathbf{x}, & \nabla h_2(\mathbf{x}) &= -\begin{pmatrix} 1 \\ 0 \end{pmatrix}, \\ d^2 f(\mathbf{x}) &= d^2 h_1(\mathbf{x}) = \mathbf{I}_2, & d^2 h_2(\mathbf{x}) &= \mathbf{0}. \end{aligned}$$

Path following starts from the unconstrained solution $\mathbf{x}(0) = \mathbf{b}$. The left panel of Fig. 16.2 plots the vector field $\frac{d}{d\rho} \mathbf{x}$ at the time $\rho = 0$. The right panel shows the solution path for projection from the points $(-2, 0.5)$,

$(-2, 1.5)$, $(-1, 2)$, $(2, 1.5)$, $(2, 0)$, $(1, 2)$, and $(-0.5, -2)$ onto the feasible region. In contrast to the previous example, small steps are taken. In projecting the point $\mathbf{b}^* = (-1, 2)$ onto $(0, 1)$, our software exploits the ODE45 solver of MATLAB. Following the solution path requires derivatives at 19 different time points. Dykstra’s algorithm by comparison takes about 30 iterations to converge. ■

16.5 Bregman Iteration

We met Bregman functions previously in Sect. 13.3 in the study of adaptive barrier methods. If $J(\mathbf{u})$ is a convex function and $\mathbf{p} \in \partial J(\mathbf{v})$ is any subgradient, then the associated Bregman function

$$D_J^{\mathbf{p}}(\mathbf{u} \mid \mathbf{v}) = J(\mathbf{u}) - J(\mathbf{v}) - \mathbf{p}^*(\mathbf{u} - \mathbf{v})$$

defines a kind of distance anchored at \mathbf{v} [24, 25]. When $J(\mathbf{u})$ is differentiable, the superscript \mathbf{p} is redundant, and we omit it. The exercises at the end of the chapter list some properties of Bregman distances. In general, the symmetry and triangle properties of a metric fail. Figure 16.3 illustrates graphically a Bregman distance for a smooth function in one dimension. Problem 9 lists some commonly encountered Bregman distances.

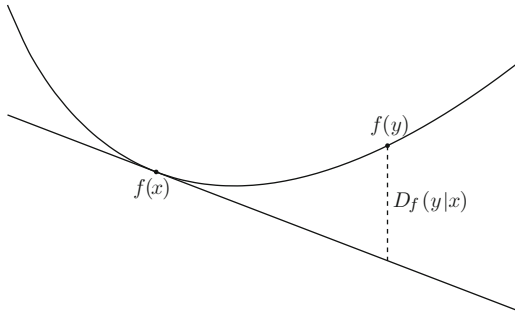


FIGURE 16.3. The Bregman distance generated by a smooth function $f(\mathbf{x})$

Example 16.5.1 Kullback–Leibler Divergence

The density function of a random vector \mathbf{Y} in a natural exponential family can be written as

$$f(\mathbf{y} \mid \boldsymbol{\theta}) = g(\mathbf{y})e^{-h(\boldsymbol{\theta})+\mathbf{y}^*\boldsymbol{\theta}}$$

for a parameter vector $\boldsymbol{\theta}$. Here \mathbf{Y} itself is the sufficient statistic. The five most important univariate examples are: the normal distribution with known variance, the Poisson distribution, the gamma distribution with

known shape parameter, the binomial distribution with known number of trials, and the negative binomial with known number of required successes. The analysis of Sect. 10.6 shows that

$$E_{\theta}(\mathbf{Y}) = \nabla h(\theta), \quad \text{Var}_{\theta}(\mathbf{Y}) = d^2 h(\theta).$$

It follows that $h(\theta)$ is convex. The Fenchel conjugate

$$h^*(\mathbf{y}) = \sup_{\theta} [\mathbf{y}^* \theta - h(\theta)] = \sup_{\theta} \ln f(\mathbf{y} | \theta)$$

determines the maximum likelihood estimate $\hat{\theta}$ through the likelihood equation $\mathbf{y} = \nabla h(\theta)$. The Bregman distance

$$\begin{aligned} D_h(\theta_1, \theta_0) &= h(\theta_1) - h(\theta_0) - dh(\theta_0)(\theta_1 - \theta_0) \\ &= E_{\theta_0} \left[h(\theta_1) - h(\theta_0) - \mathbf{Y}^*(\theta_1 - \theta_0) \right] \\ &= E_{\theta_0} \left[\ln \frac{f(\mathbf{Y} | \theta_0)}{f(\mathbf{Y} | \theta_1)} \right] \end{aligned}$$

coincides with the Kullback–Leibler divergence of the densities $f(\mathbf{y} | \theta_0)$ and $f(\mathbf{y} | \theta_1)$. ■

Osher [120, 208, 279] and colleagues have pioneered the application of Bregman iteration in compressed sensing. One of their motivating examples is to minimize the convex function $J(\mathbf{u}) = \|\mathbf{u}\|_1$ subject to the equality constraint $H(\mathbf{u}) = \min_{\mathbf{v}} H(\mathbf{v})$ for $H(\mathbf{u})$ smooth and convex. In the simple case of basis pursuit, $H(\mathbf{u})$ equals the sum of squares criterion $\frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{f}\|^2$. The next Bregman iterate \mathbf{u}_{k+1} minimizes the surrogate function

$$G(\mathbf{u} | \mathbf{u}_k) = \lambda H(\mathbf{u}) + D_J^{\mathbf{p}_k}(\mathbf{u} | \mathbf{u}_k).$$

Here $\lambda > 0$ is a scaling constant determining the relative contribution of $H(\mathbf{u})$. For the sake of simplicity, we will take $\lambda = 1$ by absorbing its value in the definition of $H(\mathbf{u})$. We will also shift $H(\mathbf{u})$ so that its minimum value is 0. This action does not affect the choice of the update \mathbf{u}_{k+1} . Bregman iteration tries to drive $H(\mathbf{u})$ to 0 and simultaneously minimize $J(\mathbf{u})$ subject to the constraint $H(\mathbf{u}) = 0$. There are four obvious questions raised by Bregman iteration. First, does the next iterate \mathbf{u}_{k+1} exist? This is certainly the case when $G(\mathbf{u} | \mathbf{u}_k)$ is coercive. Second, is \mathbf{u}_{k+1} uniquely determined? Strict convexity of $G(\mathbf{u} | \mathbf{u}_k)$ suffices. Third, how can one find \mathbf{u}_{k+1} ? In sparse problems with $J(\mathbf{u}) = \|\mathbf{u}\|_1$, coordinate descent works well. The fourth question of how to choose the subgradient \mathbf{p}_k is the most subtle of all.

The MM principle is the key to understanding Bregman iteration. Both the objective function $H(\mathbf{u})$ and the Bregman function $D_J^{\mathbf{p}_k}(\mathbf{u} | \mathbf{u}_k)$ are convex. Because $D_J^{\mathbf{p}_k}(\mathbf{u} | \mathbf{u}_k)$ is anchored at \mathbf{u}_k and majorizes the constant

0, the surrogate function $G(\mathbf{u} \mid \mathbf{u}_k)$ majorizes $H(\mathbf{u})$ at \mathbf{u}_k . Minimizing the surrogate therefore drives $H(\mathbf{u})$ downhill. The preferred update of the subgradient is straightforward to explain in this context. The convex stationarity condition $\mathbf{0} \in \nabla H(\mathbf{u}_{k+1}) + \partial J(\mathbf{u}_{k+1}) - \mathbf{p}_k$ shows that

$$\mathbf{p}_{k+1} = \mathbf{p}_k - \nabla H(\mathbf{u}_{k+1}) \in \partial J(\mathbf{u}_{k+1}).$$

Thus, the MM update furnishes a candidate subgradient. Treating the last equation recursively gives the further identity

$$\mathbf{p}_k = \mathbf{p}_0 - \sum_{i=1}^k \nabla H(\mathbf{u}_i). \tag{16.14}$$

If $J(\mathbf{u}) = \|\mathbf{u}\|_1$, then $\mathbf{u}_0 = \mathbf{p}_0 = \mathbf{0}$ is an obvious starting point for Bregman iteration.

Fortunately, there is a simple proof that the Bregman iterates as just defined send $H(\mathbf{u})$ to 0. The argument invokes the identity

$$D_J^p(\mathbf{u} \mid \mathbf{v}) + D_J^q(\mathbf{v} \mid \mathbf{w}) - D_J^q(\mathbf{u} \mid \mathbf{w}) = (\mathbf{p} - \mathbf{q})^*(\mathbf{v} - \mathbf{u}), \tag{16.15}$$

which the reader can readily verify. If we suppose that $H(\mathbf{u})$ is coercive, then it achieves its minimum of 0 at some point $\hat{\mathbf{u}}$. The identity (16.15) and the convexity of $H(\mathbf{u})$ therefore imply that

$$\begin{aligned} & D_J^{\mathbf{p}_k}(\hat{\mathbf{u}} \mid \mathbf{u}_k) - D_J^{\mathbf{p}_{k-1}}(\hat{\mathbf{u}} \mid \mathbf{u}_{k-1}) \\ & \leq D_J^{\mathbf{p}_k}(\hat{\mathbf{u}} \mid \mathbf{u}_k) + D_J^{\mathbf{p}_{k-1}}(\mathbf{u}_k \mid \mathbf{u}_{k-1}) - D_J^{\mathbf{p}_{k-1}}(\hat{\mathbf{u}} \mid \mathbf{u}_{k-1}) \\ & = (\mathbf{p}_k - \mathbf{p}_{k-1})^*(\mathbf{u}_k - \hat{\mathbf{u}}) \\ & = dH(\mathbf{u}_k)(\hat{\mathbf{u}} - \mathbf{u}_k) \\ & \leq H(\hat{\mathbf{u}}) - H(\mathbf{u}_k) \\ & = -H(\mathbf{u}_k). \end{aligned} \tag{16.16}$$

Summing the extremes of inequality (16.16) from 1 to m produces

$$\begin{aligned} \sum_{k=1}^m [D_J^{\mathbf{p}_k}(\hat{\mathbf{u}} \mid \mathbf{u}_k) - D_J^{\mathbf{p}_{k-1}}(\hat{\mathbf{u}} \mid \mathbf{u}_{k-1})] &= D_J^{\mathbf{p}_m}(\hat{\mathbf{u}} \mid \mathbf{u}_m) - D_J^{\mathbf{p}_0}(\hat{\mathbf{u}} \mid \mathbf{u}_0) \\ &\leq -\sum_{k=1}^m H(\mathbf{u}_k) \\ &\leq -mH(\mathbf{u}_m). \end{aligned}$$

Rearranging this now yields

$$H(\mathbf{u}_m) \leq \frac{1}{m} [D_J^{\mathbf{p}_0}(\hat{\mathbf{u}} \mid \mathbf{u}_0) - D_J^{\mathbf{p}_m}(\hat{\mathbf{u}} \mid \mathbf{u}_m)] \leq \frac{1}{m} D_J^{\mathbf{p}_0}(\hat{\mathbf{u}} \mid \mathbf{u}_0).$$

The convergence of $H(\mathbf{u}_m)$ to 0 is an immediate consequence.

For the sum of squares criterion $H(\mathbf{u}) = \frac{\lambda}{2} \|\mathbf{A}\mathbf{u} - \mathbf{f}\|^2$, Bregman iteration simplifies considerably. Given the initial conditions $\mathbf{u}_0 = \mathbf{p}_0 = \mathbf{0}$, define $\mathbf{f}_0 = \mathbf{f}$ and $\mathbf{f}_k = \mathbf{f}_{k-1} + (\mathbf{f} - \mathbf{A}\mathbf{u}_k)$. It then follows from equation (16.14) and telescoping that

$$\begin{aligned} H(\mathbf{u}) - \mathbf{p}_k^* \mathbf{u} &= H(\mathbf{u}) + \left[\sum_{i=1}^k \nabla H(\mathbf{u}_i) \right]^* \mathbf{u} \\ &= \frac{\lambda}{2} \|\mathbf{A}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{f}\|^2 - \lambda \mathbf{f}^* \mathbf{A}\mathbf{u} + \lambda \left[\sum_{i=1}^k \mathbf{A}^* (\mathbf{A}\mathbf{u}_i - \mathbf{f}) \right]^* \mathbf{u} \\ &= \frac{\lambda}{2} \|\mathbf{A}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{f}\|^2 - \lambda \left[\mathbf{f} + \sum_{i=1}^k (\mathbf{f} - \mathbf{A}\mathbf{u}_i) \right]^* \mathbf{A}\mathbf{u} \\ &= \frac{\lambda}{2} \|\mathbf{A}\mathbf{u}\|^2 + \frac{\lambda}{2} \|\mathbf{f}\|^2 - \lambda \mathbf{f}_k^* \mathbf{A}\mathbf{u} \\ &= \frac{\lambda}{2} \|\mathbf{A}\mathbf{u} - \mathbf{f}_k\|^2 - \frac{\lambda}{2} \|\mathbf{f}_k\|^2 + \frac{\lambda}{2} \|\mathbf{f}\|^2. \end{aligned}$$

Thus, the Bregman surrogate function becomes

$$G(\mathbf{u} \mid \mathbf{u}_k) = \frac{\lambda}{2} \|\mathbf{A}\mathbf{u} - \mathbf{f}_k\|^2 + J(\mathbf{u}) + c_k,$$

where c_k is an irrelevant constant. Section 13.5 sketches how to find \mathbf{u}_{k+1} by coordinate descent when $J(\mathbf{u}) = \|\mathbf{u}\|_1$.

The linearized version of Bregman iteration [209] relies on the approximation

$$H(\mathbf{u}) \approx H(\mathbf{u}_k) + dH(\mathbf{u}_k)(\mathbf{u} - \mathbf{u}_k).$$

Since this Taylor expansion is apt to be accurate only when $\|\mathbf{u} - \mathbf{u}_k\|$ is small, the surrogate function is re-defined as

$$G(\mathbf{u} \mid \mathbf{u}_k) = H(\mathbf{u}_k) + dH(\mathbf{u}_k)(\mathbf{u} - \mathbf{u}_k) + D_{J^k}^{\mathbf{p}}(\mathbf{u} \mid \mathbf{u}_k) + \frac{1}{2\delta} \|\mathbf{u} - \mathbf{u}_k\|^2$$

for $\delta > 0$. The quadratic penalty majorizes 0 and shrinks the next iterate \mathbf{u}_{k+1} toward \mathbf{u}_k . Separation of parameters in the surrogate function when $J(\mathbf{u}) = \|\mathbf{u}\|_1$ is a huge advantage in solving for \mathbf{u}_{k+1} . Examination of the convex stationarity conditions then shows that

$$u_{k+1,j} = \begin{cases} u_{kj}^+ = u_{kj} + \delta \left[p_{kj} - \frac{\partial}{\partial u_j} H(\mathbf{u}_k) - 1 \right] & u_{kj}^+ > 0 \\ u_{kj}^- = u_{kj} + \delta \left[p_{kj} - \frac{\partial}{\partial u_j} H(\mathbf{u}_k) + 1 \right] & u_{kj}^- < 0 \\ 0 & \text{otherwise.} \end{cases}$$

Given that the surrogate function approximately majorizes $H(\mathbf{u})$ in a neighborhood of \mathbf{u}_k , one can also confidently expect the linearized Bregman iterates to drive $H(\mathbf{u})$ downhill.

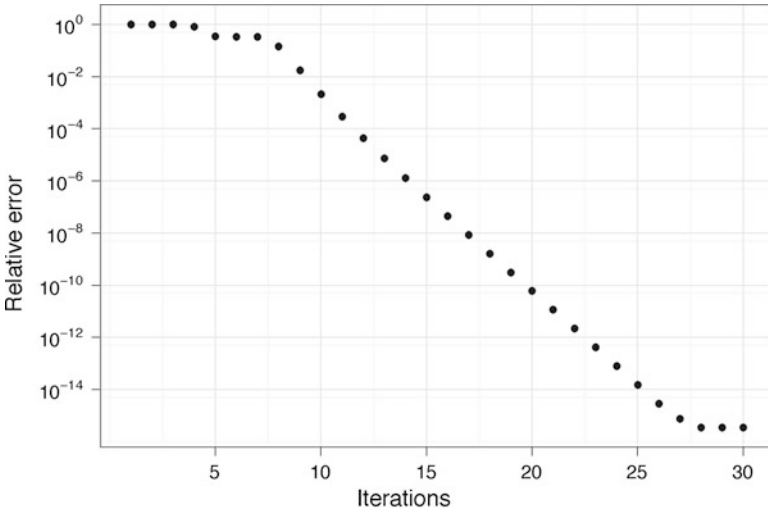


FIGURE 16.4. Relative error versus iteration number for basis pursuit

Figure 16.4 plots the results of linearized Bregman iteration for a simple numerical example. Here the entries of the $10^2 \times 10^5$ matrix \mathbf{A} are populated with independent standard normal deviates. The generating sparse vector \mathbf{u} has all entries equal to 0 except for

$$u_{45373} = -1.162589, \quad u_{57442} = 2.436616, \quad u_{81515} = 1.876241.$$

The response vector \mathbf{f} equals $\mathbf{A}\mathbf{u}$, and the Bregman constants are $\delta = 0.01$ and $\lambda = 1000$. The figure plots the relative error $\|\mathbf{u} - \mathbf{u}_k\| / \|\mathbf{u}\|$ as a function of iteration number k . It is remarkable how quickly the true \mathbf{u} is recovered in the absence of noise. As one might expect, lasso penalized linear regression, also known as basis pursuit denoising, produces solutions very similar to basis pursuit.

16.6 Split Bregman Iteration

Our focus on the penalty $J(\mathbf{u}) = \|\mathbf{u}\|_1$ obscures the fact that many applications really require $J(\mathbf{u}) = \|\mathbf{D}\mathbf{u}\|_1$ for a constant matrix \mathbf{D} . For instance, the well-known image denoising model of Rudin, Osher, and Fatemi [224] minimizes the total variation regularized sum of squares criterion

$$\frac{1}{2} \sum_{i,j} (w_{ij} - u_{ij})^2 + \rho \sum_{i,j} \sqrt{(u_{i+1,j} - u_{ij})^2 + (u_{i,j+1} - u_{ij})^2},$$

where w_{ij} is the corrupted intensity and u_{ij} is the true intensity for pixel (i, j) of an image. The total variation penalty represented by the second

sum is intended to smooth the reconstructed image while preserving its edges. A similar effect can be achieved by adopting the anisotropic penalty

$$\rho \sum_{i,j} \left(|u_{i+1,j} - u_{ij}| + |u_{i,j+1} - u_{ij}| \right),$$

which has the form $J(\mathbf{u}) = \rho \|\mathbf{D}\mathbf{u}\|_1$ for $\mathbf{D}\mathbf{u}$ linear.

Split Bregman iteration is intended to handle this kind of situation [106]. Consider minimization of the criterion $E(\mathbf{u}) + \|\mathbf{D}\mathbf{u}\|_1$, where \mathbf{D} is linear and $E(\mathbf{u})$ is convex and differentiable. In split Bregman iteration the general idea is to introduce a new variable $\mathbf{d} = \mathbf{D}\mathbf{u}$ and carry out Bregman iteration with the objective function $H(\mathbf{u}, \mathbf{d}) = \frac{1}{2\delta} \|\mathbf{d} - \mathbf{D}\mathbf{u}\|^2$ modified by the Bregman function derived from $J(\mathbf{u}, \mathbf{d}) = E(\mathbf{u}) + \|\mathbf{d}\|_1$. Thus, one selects the pair $(\mathbf{u}_{k+1}, \mathbf{d}_{k+1})$ to minimize the criterion

$$\frac{1}{2\delta} \|\mathbf{d} - \mathbf{D}\mathbf{u}\|^2 + E(\mathbf{u}) + \|\mathbf{d}\|_1 - \mathbf{p}_k^*(\mathbf{u} - \mathbf{u}_k) - \mathbf{q}_k^*(\mathbf{d} - \mathbf{d}_k) \quad (16.17)$$

for subgradients $\mathbf{p}_k \in \partial E(\mathbf{u}_k)$ and $\mathbf{q}_k \in \partial \|\mathbf{d}_k\|_1$. Once the next iterate $(\mathbf{u}_{k+1}, \mathbf{d}_{k+1})$ is determined, the new subgradients

$$\begin{aligned} \mathbf{p}_{k+1} &= \mathbf{p}_k - \nabla_{\mathbf{u}} H(\mathbf{u}_{k+1}, \mathbf{d}_{k+1}) \\ \mathbf{q}_{k+1} &= \mathbf{q}_k - \nabla_{\mathbf{d}} H(\mathbf{u}_{k+1}, \mathbf{d}_{k+1}) \end{aligned}$$

are defined. Block relaxation is the natural method of minimizing the criterion (16.17). If $E(\mathbf{u})$ is well behaved, then one can minimize

$$\frac{1}{2\delta} \|\mathbf{d} - \mathbf{D}\mathbf{u}\|^2 + E(\mathbf{u}) - \mathbf{p}_k^*(\mathbf{u} - \mathbf{u}_k)$$

with respect to \mathbf{u} by Newton's method. Minimizing

$$\frac{1}{2\delta} \|\mathbf{d} - \mathbf{D}\mathbf{u}\|^2 + \|\mathbf{d}\|_1 - \mathbf{q}_k^*(\mathbf{d} - \mathbf{d}_k)$$

with respect to \mathbf{d} can be achieved in a single iteration by the shrinkage rule

$$d_j = \begin{cases} d_j^+ = (\mathbf{D}\mathbf{u})_j + \delta(q_{kj} - 1) & d_j^+ > 0 \\ d_j^- = (\mathbf{D}\mathbf{u})_j + \delta(q_{kj} + 1) & d_j^- < 0 \\ 0 & \text{otherwise} \end{cases}$$

suggested in our discussion of linearized Bregman iteration.

Example 16.6.1 The Fused Lasso

The fused lasso [257] penalty is the ℓ_1 norm of the successive differences between parameters. The simplest possible fused lasso problem minimizes the penalized sum of squares criterion

$$E(\mathbf{u}) + \|\mathbf{D}\mathbf{u}\|_1 = \frac{\lambda}{2} \sum_{i=1}^n (y_i - u_i)^2 + \sum_{i=1}^{n-1} |u_{i+1} - u_i|. \quad (16.18)$$

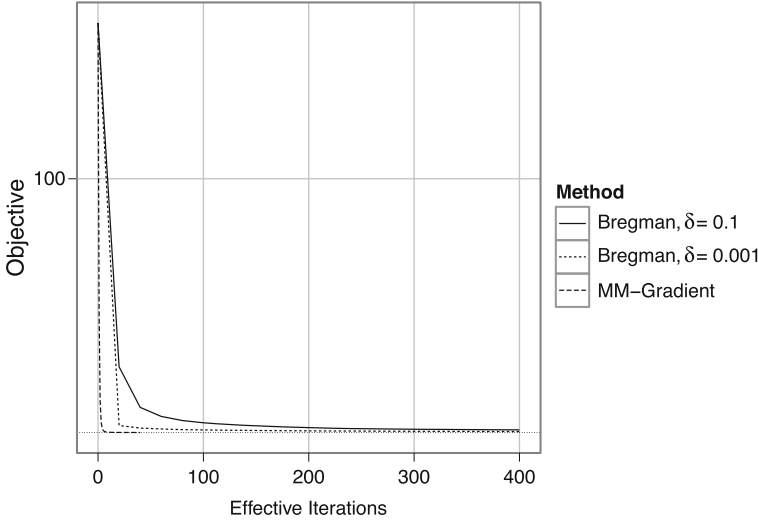


FIGURE 16.5. Convergence to the objective (16.18) for the fused lasso problem

Note that the matrix \mathbf{D} in the penalty $\|\mathbf{D}\mathbf{u}\|_1$ is bidiagonal. The differences $u_{i+1} - u_i$ under the absolute value signs make it difficult to implement coordinate descent, so split Bregman iteration is an attractive possibility. An alternative [280] is to change the penalty slightly and minimize the revised objective function

$$f(\mathbf{u}) = \frac{\lambda}{2} \sum_{i=1}^n (y_i - u_i)^2 + \sum_{i=1}^{n-1} [(u_{i+1} - u_i)^2 + \epsilon]^{1/2},$$

which smoothly approximates the original objective function for small positive ϵ . The majorization (8.12) translates into the quadratic majorization

$$g(\mathbf{u} \mid \mathbf{u}_k) = \frac{\lambda}{2} \sum_{i=1}^n (y_i - u_i)^2 + \sum_{i=1}^{n-1} \frac{(u_{i+1} - u_i)^2}{2[u_{k,i+1} - u_{k,i}]^2 + \epsilon]^{1/2}} + c$$

for an irrelevant constant c . The MM gradient algorithm for updating \mathbf{u} is easy to implement because the second differential $d^2g(\mathbf{u} \mid \mathbf{u}_k)$ is tridiagonal and efficiently inverted at $\mathbf{u} = \mathbf{u}_k$ by Thomas’s algorithm [52]. In fact, one step of Newton’s method minimizes the quadratic $g(\mathbf{u} \mid \mathbf{u}_k)$. Updating \mathbf{u} in split Bregman iteration also benefits from Thomas’s algorithm. Observe, however, that the multiple inner iterations of block descent puts split Bregman iteration at a computational disadvantage.

Figure 16.5 compares the performance of split Bregman iteration and the MM gradient algorithm on simulated data with independently generated responses y_1, \dots, y_{1000} sampled from two Gaussian densities. These densities share a common standard deviation of 0.2 but differ in their means of

0.6 for $501 \leq i \leq 550$ and 0 otherwise. Thus, there is a brief elevation in the signal for a short segment in the middle. Both algorithms commence from the starting values $\mathbf{u}_0 = \mathbf{y}$ and $\mathbf{d}_0 = \mathbf{0}$. For split Bregman iteration $\mathbf{p}_0 = \mathbf{q}_0 = \mathbf{0}$; for the MM gradient algorithm $\epsilon = 10^{-10}$. The tuning constant $\lambda = 2.5/\sqrt{\ln 1000}$. The effective number of iterations plotted in Fig. 16.5 counts the total inner iterations for split Bregman iteration. ■

16.7 Convergence of Bregman Iteration

Proving convergence of Bregman iteration is difficult [147, 279]. The crux of the problem is that minimizing $J(\mathbf{u})$ is secondary to minimizing $H(\mathbf{u})$. Owing to the difficulties, we will only tackle convergence for the special choices $H(\mathbf{u}) = \frac{\lambda}{2}\|\mathbf{A}\mathbf{u} - \mathbf{f}\|^2$ and $J(\mathbf{u}) = \|\mathbf{D}\mathbf{u}\|_1$ made in image denoising. If \mathbf{A} does not have full column rank, then the purpose of Bregman iteration is to minimize the secondary criterion $J(\mathbf{u}) = \|\mathbf{D}\mathbf{u}\|_1$ subject to $\mathbf{A}\mathbf{u} = \mathbf{f}$. The first question that comes to mind is whether the minimum exists. Let K be the kernel of the matrix \mathbf{A} and \mathbf{x} be a particular solution of the equation $\mathbf{A}\mathbf{u} = \mathbf{f}$. This equation's solution space is simply $\mathbf{x} + K$. If $\mathbf{D}\mathbf{x} = \mathbf{0}$, then the minimum value of $J(\mathbf{u})$ is achieved at \mathbf{x} . If in addition $\mathbf{D}\mathbf{y} = \mathbf{0}$ for some $\mathbf{y} \in K$, then the minimum of 0 is also achieved along the entire line through \mathbf{x} along the direction \mathbf{y} .

Now consider whether the function $\mathbf{y} \mapsto \|\mathbf{D}(\mathbf{x} + \mathbf{y})\|_1$ is coercive on K . Coerciveness implies that the minimum exists. In view of the equivalence of norms, it suffices to decide whether the function $\mathbf{y} \mapsto \|\mathbf{D}(\mathbf{x} + \mathbf{y})\|$ is coercive. Proposition 12.3.1 implies that the Euclidean norm is coercive if and only if

$$\lim_{t \rightarrow \infty} \|\mathbf{D}(\mathbf{x} + t\mathbf{y})\|^2 = \lim_{t \rightarrow \infty} \|\mathbf{D}\mathbf{x}\|^2 + 2t\mathbf{y}^* \mathbf{D}^* \mathbf{D} \mathbf{x} + t^2 \mathbf{y}^* \mathbf{D}^* \mathbf{D} \mathbf{y} = \infty$$

for all $\mathbf{y} \in K$. This is clearly equivalent to the condition $\mathbf{D}\mathbf{y} \neq \mathbf{0}$ for all $\mathbf{y} \in K$. Hence, the condition $\mathbf{D}\mathbf{y} \neq \mathbf{0}$ for every $\mathbf{y} \in K$ is necessary and sufficient for coerciveness. According to the discussion following Example 5.5.3, another equivalent condition is that the matrix $\mathbf{A}^* \mathbf{A} + \gamma \mathbf{D}^* \mathbf{D}$ is positive definite for some $\gamma > 0$.

The boundedness of the iterates and subdifferentials plays a key role in proving convergence. The latter are easy to handle because the chain rule entails $\partial J(\mathbf{u}) = \mathbf{D}^* \partial \|\mathbf{D}\mathbf{u}\|_1$. For any $\mathbf{v} \in \mathbb{R}^n$, $\partial \|\mathbf{v}\|_1$ is contained in the n -dimensional cube $[-1, 1]^n$. Hence, $\partial J(\mathbf{u})$ is contained in the compact set $\mathbf{D}^*[-1, 1]^n$ for all \mathbf{u} . One can guarantee that the iterates \mathbf{u}_k are also bounded whenever $H(\mathbf{u})$ is coercive. Unfortunately, this is not the case for an under-determined system $\mathbf{A}\mathbf{u} = \mathbf{f}$. We will simply postulate that the iteration sequence \mathbf{u}_k is bounded. We will also assume that $J(\mathbf{u})$ achieves its constrained minimum at some point \mathbf{w} .

Now suppose some subsequence of the Bregman iterates \mathbf{u}_{k+1} converges to a cluster point \mathbf{y} . To show that \mathbf{y} minimizes the criterion $J(\mathbf{u})$ subject to $H(\mathbf{u}) = 0$, we examine the inequality

$$\begin{aligned} H(\mathbf{u}_{k+1}) + J(\mathbf{u}_{k+1}) - J(\mathbf{u}_k) - \mathbf{p}_k^*(\mathbf{u}_{k+1} - \mathbf{u}_k) \\ \leq H(\mathbf{w}) + J(\mathbf{w}) - J(\mathbf{u}_k) - \mathbf{p}_k^*(\mathbf{w} - \mathbf{u}_k). \end{aligned}$$

By passing to a subsequence if necessary, we can assume that \mathbf{u}_k converges to \mathbf{x} , \mathbf{u}_{k+1} converges to \mathbf{y} , and \mathbf{p}_k converges to \mathbf{p} . If we can show that

$$\mathbf{p}^*(\mathbf{y} - \mathbf{x}) = \mathbf{p}^*(\mathbf{w} - \mathbf{x}) = 0,$$

then in the limit the descent property of $H(\mathbf{u})$ implies $J(\mathbf{y}) \leq J(\mathbf{w})$. Thus, the cluster point \mathbf{y} is also optimal.

We now check that $\mathbf{p}^*(\mathbf{y} - \mathbf{x}) = 0$. The other inner product vanishes for similar reasons. Recall that we commence with $\mathbf{u}_0 = \mathbf{p}_0 = \mathbf{0}$. Thus, \mathbf{p}_0 belongs to the range of \mathbf{A}^* . In general, this assertion is true for all \mathbf{p}_k because $\mathbf{p}_{k+1} = \mathbf{p}_k - \lambda \mathbf{A}^*(\mathbf{A}\mathbf{u}_{k+1} - \mathbf{f})$. Given that the range of \mathbf{A}^* is closed, there exists a vector \mathbf{z} with $\mathbf{p} = \mathbf{A}^*\mathbf{z}$. Furthermore, $\mathbf{A}\mathbf{y} = \mathbf{A}\mathbf{x} = \mathbf{f}$ since $H(\mathbf{y}) = H(\mathbf{x}) = 0$. It follows that

$$\mathbf{p}^*(\mathbf{y} - \mathbf{x}) = \mathbf{z}^*\mathbf{A}(\mathbf{y} - \mathbf{x}) = 0.$$

Let us summarize this discussion in a proposition.

Proposition 16.7.1 *Consider Bregman iteration for the function choices $H(\mathbf{u}) = \frac{\lambda}{2}\|\mathbf{A}\mathbf{u} - \mathbf{f}\|^2$ and $J(\mathbf{u}) = \|\mathbf{D}\mathbf{u}\|_1$. The minimum value of $J(\mathbf{u})$ subject to $H(\mathbf{u}) = 0$ is attained provided $\mathbf{D}\mathbf{u} \neq \mathbf{0}$ whenever $\mathbf{A}\mathbf{u} = \mathbf{0}$. When the minimum is attained and the Bregman iterates \mathbf{u}_k remain bounded, every cluster point of the iterates \mathbf{u}_k achieves the minimum.*

Proof: See the foregoing comments. ■

For the basis pursuit problem, it is known that Bregman iteration converges in a finite number of steps [279]. Problems 15 and 16 sketch a proof of this important fact.

16.8 Problems

1. Suppose that the real-valued function $f(\mathbf{x})$ is twice differentiable and that $b = \sup_{\mathbf{z}} \|d^2f(\mathbf{z})\|$ is finite. Prove the Lipschitz inequality

$$\|\nabla f(\mathbf{y}) - \nabla f(\mathbf{x})\| \leq b\|\mathbf{y} - \mathbf{x}\|$$

for all \mathbf{x} and \mathbf{y} . If $f(\mathbf{x})$ satisfies the Lipschitz inequality, then prove that

$$g(\mathbf{x} \mid \mathbf{x}_m) = f(\mathbf{x}_m) + df(\mathbf{x}_m)(\mathbf{x} - \mathbf{x}_m) + \frac{b}{2}\|\mathbf{x} - \mathbf{x}_m\|^2$$

majorizes $f(\mathbf{x})$. (Hint: Expand $f(\mathbf{x})$ to first order around \mathbf{x}_m . Rearrange the integral remainder and bound.)

2. As described in Problem 1, assume the gradient of $f(\mathbf{x})$ is Lipschitz. Consider the projected gradient algorithm

$$\mathbf{x}_{m+1} = P_S \left[\mathbf{x}_m - \frac{\rho}{b} \nabla f(\mathbf{x}_m) \right]$$

for $\rho \in (0, 2)$. From the obtuse angle criterion deduce

$$df(\mathbf{x}_m)(\mathbf{x}_{m+1} - \mathbf{x}_m) \leq -\frac{b}{\rho} \|\mathbf{x}_{m+1} - \mathbf{x}_m\|^2.$$

Add this inequality to the majorizing inequality and further deduce

$$f(\mathbf{x}_{m+1}) \leq f(\mathbf{x}_m) - \left[\frac{b}{\rho} - \frac{b}{2} \right] \|\mathbf{x}_{m+1} - \mathbf{x}_m\|^2.$$

It follows that the sequence $f(\mathbf{x}_m)$ is decreasing. If $f(\mathbf{x})$ is coercive or S is compact, then argue that $\lim_{m \rightarrow \infty} f(\mathbf{x}_m)$ exists and that $\lim_{m \rightarrow \infty} \|\mathbf{x}_{m+1} - \mathbf{x}_m\| = 0$. If we suppose \mathbf{y} is a cluster point of the sequence \mathbf{x}_m , then the second of these limits shows that

$$P_S \left[\mathbf{y} - \frac{\rho}{b} \nabla f(\mathbf{y}) \right] = \mathbf{y}.$$

Apply the obtuse angle criterion to any $\mathbf{z} \in S$, and deduce the necessary condition $df(\mathbf{y})(\mathbf{z} - \mathbf{y}) \geq 0$ for optimality. When $f(\mathbf{x})$ is also convex, conclude that \mathbf{y} provides a global minimum. If $f(\mathbf{x})$ is strictly convex, then $\lim_{m \rightarrow \infty} \mathbf{x}_m$ exists and furnishes the unique global minimum point.

3. Prove that the function $\mathcal{E}_\rho(\mathbf{y})$ appearing in the exact penalty method is convex whenever $f(\mathbf{y})$ and the inequality constraints $h_j(\mathbf{y})$ are convex and the equality constraints $g_i(\mathbf{y})$ are affine.
4. In the exact penalty method for projection onto the half-disc, show that the solution path initially: (a) heads toward the origin when $\mathbf{x} \in \{\mathbf{x} : \|\mathbf{x}\|^2 > 1, x_1 > 0\}$ or $\mathbf{x} \in \{\mathbf{x} : |x_2| > 1, x_1 = 0\}$, (b) heads toward the point $(1, 0)^*$ when $\mathbf{x} \in \{\mathbf{x} : \|\mathbf{x}\|^2 > 1, x_1 < 0\}$, (c) follows the unit circle when $\mathbf{x} \in \{\mathbf{x} : \|\mathbf{x}\|^2 = 1, x_1 < 0\}$, and (d) heads toward the x_2 -axis when $\mathbf{x} \in \{\mathbf{x} : \|\mathbf{x}\|^2 < 1, x_1 < 0\}$. (Hint: See the left panel of Fig. 16.2.)
5. Path following can be conducted without invoking the exact penalty method [31, 46]. Consider minimizing the convex differentiable function $f(\mathbf{x})$ subject to the standard linear programming constraints $\mathbf{x} \geq \mathbf{0}$ and $\mathbf{A}\mathbf{x} = \mathbf{b}$. Given an initial feasible point $\mathbf{x}_0 > \mathbf{0}$, one can devise a differential equation $\frac{d}{dt}x(t) = G(\mathbf{x})$ whose solution $x(t)$ is likely

to converge to the optimal point. Simply take $G(\mathbf{x}) = D(\mathbf{x})P(\mathbf{x})v(\mathbf{x})$, where $D(\mathbf{x}) = \text{diag}(\mathbf{x})$ is a diagonal matrix with diagonal entries given by the vector \mathbf{x} and $P(\mathbf{x})$ is orthogonal projection onto the null space of $\mathbf{A}D(\mathbf{x})$. The matrix $D(\mathbf{x})$ slows the trajectory down as it approaches a boundary $x_i = 0$. The matrix $P(\mathbf{x})$ ensures that the value of $\mathbf{A}x(t)$ remains fixed at the constant \mathbf{b} . Check this fact. Show that the choice $v(\mathbf{x}) = -Q(\mathbf{x})P(\mathbf{x})D(\mathbf{x})\nabla f(\mathbf{x})$ for $Q(\mathbf{x})$ positive semidefinite yields $\frac{d}{dt}f[x(t)] \leq 0$. In other words, $f(\mathbf{x})$ is a Liapunov function for the solution path. For an analogue of steepest descent useful in linear programming, the choice $Q(\mathbf{x}) = \mathbf{I}$ is obvious. For an analogue of Newton's method, the choice $Q(\mathbf{x}) = d^2f(\mathbf{x})^{-1}$ has merit. In a statistical setting where $-f(\mathbf{x})$ is a loglikelihood, \mathbf{x} is a parameter vector, and $-\nabla f(\mathbf{x})$ is the score, substitution of the expected information matrix for the observed information matrix is also reasonable.

6. Implement the path following algorithm of Problem 5 for linear programming, linear regression, or linear logistic regression. You may use the differential equation solver of Matlab or Euler's method. Cheney [46] mentions some tactics for linear programming that ease the computational burden and make the overall algorithm more stable.
7. Homotopy methods follow solution paths. Suppose you want to solve the equation $f(\mathbf{x}) = \mathbf{0}$. Choose any point \mathbf{x}_0 and define the homotopy

$$h(t, \mathbf{x}) = tf(\mathbf{x}) + (1-t)[f(\mathbf{x}) - f(\mathbf{x}_0)] = f(\mathbf{x}) + (t-1)f(\mathbf{x}_0).$$

Note that $h(0, \mathbf{x}) = f(\mathbf{x}) - f(\mathbf{x}_0)$ and $h(1, \mathbf{x}) = f(\mathbf{x})$. Furthermore, $h(0, \mathbf{x}_0) = \mathbf{0}$. Path following starts at time $t = 0$ and $\mathbf{x} = \mathbf{x}_0$. Apply the implicit function theorem to the equation $h(t, \mathbf{x}) = \mathbf{0}$, and demonstrate that a continuously differentiable solution path $x(t)$ exists satisfying the differential equation

$$\frac{d}{dt}x(t) = -df[x(t)]^{-1}f(\mathbf{x}_0). \quad (16.19)$$

8. Continuing Problem 7, suppose

$$f(\mathbf{x}) = \begin{pmatrix} \sin x_1 + e^{x_2} - 3 \\ (x_2 + 3)^2 - x_1 - 4 \end{pmatrix}.$$

Numerically integrate the differential equation (16.19) starting at the point $\mathbf{x}_0 = (5, 3)^*$. You should have an approximate zero at the time $t = 1$. If necessary, polish this approximate solution by Newton's method. The purpose of path following is to reliably reach a neighborhood of a solution.

9. Show that the four functions $J_1(\mathbf{u}) = \|\mathbf{u}\|^2$, $J_2(\mathbf{u}) = -\sum_i \log u_i$, $J_3(\mathbf{u}) = \sum_i u_i \ln u_i$, and $J_4(\mathbf{u} \mid \mathbf{v}) = \mathbf{u}^* \mathbf{A} \mathbf{u}$ generate the Bregman distances

$$\begin{aligned} D_{J_1}(\mathbf{u} \mid \mathbf{v}) &= \|\mathbf{u} - \mathbf{v}\|^2 \\ D_{J_2}(\mathbf{u} \mid \mathbf{v}) &= \sum_i \left[\frac{u_i}{v_i} - \log \left(\frac{u_i}{v_i} \right) - 1 \right] \\ D_{J_3}(\mathbf{u} \mid \mathbf{v}) &= \sum_i u_i \ln \left(\frac{u_i}{v_i} \right) - \sum_i (u_i - v_i) \\ D_{J_4}(\mathbf{u} \mid \mathbf{v}) &= (\mathbf{u} - \mathbf{v})^* \mathbf{A} (\mathbf{u} - \mathbf{v}). \end{aligned}$$

For $J_4(\mathbf{u})$ assume \mathbf{A} is positive semidefinite.

10. Prove the generalized Pythagorean identity (16.15).
11. The Bregman function $D_J^p(\mathbf{u} \mid \mathbf{v})$ has some of the properties of a distance. Show that it is 0 when $\mathbf{u} = \mathbf{v}$ and nonnegative when $\mathbf{u} \neq \mathbf{v}$. If $J(\mathbf{u})$ is strictly convex, then show that $D_J^p(\mathbf{u} \mid \mathbf{v})$ is positive when $\mathbf{u} \neq \mathbf{v}$. Also prove that $D_J^p(\mathbf{u} \mid \mathbf{v})$ is convex in its argument \mathbf{u} . Finally, prove that $D_J^p(\mathbf{u} \mid \mathbf{v}) \geq D_J^p(\mathbf{w} \mid \mathbf{v})$ when \mathbf{w} lies on the line segment between \mathbf{u} and \mathbf{v} .
12. For differentiable functions, demonstrate the Bregman identities

$$\begin{aligned} D_{cJ}(\mathbf{u} \mid \mathbf{v}) &= cD_J(\mathbf{u} \mid \mathbf{v}) \quad \text{for } c \geq 0 \\ D_{J_1+J_2}(\mathbf{u} \mid \mathbf{v}) &= D_{J_1}(\mathbf{u} \mid \mathbf{v}) + D_{J_2}(\mathbf{u} \mid \mathbf{v}) \\ D_J(\mathbf{u} \mid \mathbf{v}) &= 0 \quad \text{for } J(\mathbf{u}) \text{ affine.} \end{aligned}$$

13. Suppose X is a random variable with mean μ and $f(x)$ is a convex function defined on \mathbb{R} . Prove Jensen's inequality in the form

$$\mathbb{E}[f(X)] - f(\mu) = \mathbb{E}[D_f^p(X \mid \mu)] \geq 0,$$

where p is any subgradient of $f(x)$ at μ .

14. Suppose the convex function $f(\mathbf{x})$ and its Fenchel conjugate $f^*(\mathbf{y})$ are both differentiable. Let $\mathbf{u} = \nabla f^*(\mathbf{u}^*)$ and $\mathbf{v}^* = \nabla f(\mathbf{v})$. Prove the duality result

$$D_f(\mathbf{u} \mid \mathbf{v}) = D_{f^*}(\mathbf{v}^* \mid \mathbf{u}^*).$$

(Hint: Recall when equality holds in the Fenchel-Young inequality.)

15. Assume the Bregman iterate \mathbf{u}_k satisfies $H(\mathbf{u}_k) = \frac{\lambda}{2} \|\mathbf{A} \mathbf{u}_k - \mathbf{f}\|^2 = 0$. Prove that \mathbf{u}_k also minimizes $J(\mathbf{u})$ subject to $H(\mathbf{u}) = 0$. (Hint: Let $\hat{\mathbf{u}}$ be optimal for $J(\mathbf{u})$ given the constraint $H(\mathbf{u}) = 0$. Note that $J(\mathbf{u}_k) \leq J(\hat{\mathbf{u}}) - \mathbf{p}_k^*(\hat{\mathbf{u}} - \mathbf{u}_k)$ and \mathbf{p}_k belongs to the range of \mathbf{A}^* .)

16. For the basis pursuit problem, let (I_+^j, I_-^j, E^j) be a partition of the index set $\{1, 2, \dots, n\}$, and define

$$\begin{aligned} U^j &= \{\mathbf{u} \in \mathbb{R}^n : u_i \geq 0, i \in I_+^j; u_i \leq 0, i \in I_-^j; u_i = 0, i \in E^j\} \\ H^j &= \inf \left\{ \frac{1}{2} \|\mathbf{A}\mathbf{u} - \mathbf{f}\|^2 : \mathbf{u} \in U^j \right\}. \end{aligned}$$

Show that there are a finite number of distinct partitions U^j and that their union equals \mathbb{R}^n . At iteration k define

$$\begin{aligned} I_+^k &= \{i : p_{ki} = 1\} \\ I_-^k &= \{i : p_{ki} = -1\} \\ E^k &= \{i : p_{ki} \in (-1, 1)\}. \end{aligned}$$

Demonstrate that $\mathbf{p}_k \in \partial J(\mathbf{u}_k)$ implies $\mathbf{u}_k \in U^k$ and that $\mathbf{u}_k \in U^j$ with $H^j > 0$ can happen for only finitely many k . Hence, for some k we have $\mathbf{u}_k \in U^j$ with $H^j = 0$. Show that this entails $\mathbf{A}\mathbf{u}_{k+1} = \mathbf{f}$. Now invoke problem 15 to prove that Bregman iteration converges in a finite number of steps.

17. In Bregman iteration suppose $H(\mathbf{u})$ achieves its minimum of 0 at the point $\hat{\mathbf{u}}$. If $H(\mathbf{u}_{k-1}) > 0$, then prove that

$$D_J^{\mathbf{p}_k}(\hat{\mathbf{u}} \mid \mathbf{u}_k) < D_J^{\mathbf{p}_{k-1}}(\hat{\mathbf{u}} \mid \mathbf{u}_{k-1}).$$

(Hint: Consider inequality (16.16).)

18. Suppose the convex functions $H(\mathbf{u})$ and $J(\mathbf{u})$ in our discussion of Bregman iteration are differentiable. Prove that the surrogate function $G(\mathbf{u} \mid \mathbf{u}_k) = H(\mathbf{u}) + D_J(\mathbf{u} \mid \mathbf{u}_k)$ satisfies

$$G(\mathbf{u} \mid \mathbf{u}_{k-1}) = G(\mathbf{u}_k \mid \mathbf{u}_{k-1}) + D_H(\mathbf{u} \mid \mathbf{u}_k) + D_J(\mathbf{u} \mid \mathbf{u}_k).$$

(Hint: $\mathbf{0} = \nabla H(\mathbf{u}_k) + \nabla J(\mathbf{u}_k) - \nabla J(\mathbf{u}_{k-1})$.)

17

The Calculus of Variations

17.1 Introduction

The calculus of variations deals with infinite dimensional optimization problems. Seventeenth century mathematicians and physicists such as Newton, Galileo, Huygens, John and James Bernoulli, L'Hôpital, and Leibniz posed and solved many variational problems. In the eighteenth century Euler made more definitive strides that were clarified and extended by Lagrange. In the nineteenth and twentieth centuries, the intellectual stimulus offered by the calculus of variations was instrumental in the development of functional analysis and control theory. Some of this rich history is explored in the books [240, 258].

The current chapter surveys the classical and more elementary parts of the calculus of variations. The subject matter is optimization of functionals such as

$$F(\mathbf{x}) = \int_a^b f[t, x(t), x'(t)] dt \quad (17.1)$$

depending on a continuously differentiable function $x(t)$ over an interval $[a, b]$. Euler and Lagrange were able to deduce that the solution satisfies the differential equation

$$\frac{d}{dt} \frac{\partial}{\partial x'} f[t, x(t), x'(t)] = \frac{\partial}{\partial x} f[t, x(t), x'(t)]. \quad (17.2)$$

If constraints are imposed on the solution, then these constraints enter the Euler-Lagrange equation via multipliers. Thus, the theory parallels the finite-dimensional case.

However, as one might expect, the theory is harder. Proof strategies based on compactness often fail while strategies based on convexity usually succeed. Much of the theory involving differentials fortunately generalizes. We tackle this theory for normed vector spaces. Most other introductions to the calculus of variations substitute a weaker version of differentiation that relies entirely on directional derivatives. In our view this forfeits the chance to bridge the gap between advanced calculus and functional analysis. Having expended so much energy on developing Carathéodory's version of the differential, we continue to pursue that definition here. Readers interested in the more traditional perspective can consult the references [46, 102, 228, 233, 240]

17.2 Classical Problems

As motivating examples, we briefly discuss some of the classical problems of the calculus of variations. Finding a solution to one of these problems is often helped by a judicious choice of a coordinate system.

Example 17.2.1 *Geodesics*

A geodesic is the shortest path between two points. In the absence of constraints, a geodesic is a straight line. Suppose the points in question are $\mathbf{p} = \mathbf{0}$ and \mathbf{q} in \mathbb{R}^n . A path is a differentiable curve \mathbf{x} from $[0, 1]$ starting at $\mathbf{0}$ and ending at \mathbf{q} . To prove that the optimal path is a straight line, we must show that $y(t) = t\mathbf{q}$ minimizes the functional

$$G(\mathbf{x}) = \int_0^1 \|\mathbf{x}'(t)\| dt.$$

But this is obvious from the inequality

$$\|\mathbf{q}\| = \left\| \int_0^1 \mathbf{x}'(t) dt \right\| \leq \int_0^1 \|\mathbf{x}'(t)\| dt.$$

A somewhat harder problem is to show that a geodesic on a sphere follows a great circle. If the sphere has radius r , then a feasible path $x(t)$ must satisfy $\|x(t)\| = r$ for all t . It is convenient to pass to spherical coordinates and assume that the sphere resides in ordinary space \mathbb{R}^3 with its center at the origin. It is also convenient to take the initial point \mathbf{p} as the north pole and parameterize the azimuthal angle $\theta(\phi)$ of a path by the polar angle ϕ , where $\phi \in [0, \pi]$ and $\theta \in [0, 2\pi]$. The path $y(\phi) = (r, \theta(\phi), \phi)^*$

in spherical coordinates automatically satisfies the radial constraint. The usual arguments from elementary calculus show that $y(\phi)$ has arclength

$$G(\mathbf{y}) = r \int_0^{\phi_1} \sqrt{1 + [\theta'(\phi) \sin \phi]^2} d\phi$$

between the north pole and $\mathbf{q} = (r, \theta_1, \phi_1)$. It is clear that $G(\mathbf{y})$ is minimized by taking $\theta(\phi)$ equal to the constant θ_1 . This implies that the solution is an arc of a great circle. ■

Example 17.2.2 *Minimal Surface Area of Revolution*

In the plane \mathbb{R}^2 , imagine rotating a curve $y(x)$ about the x -axis. This generates a surface of revolution with area

$$S(\mathbf{y}) = 2\pi \int_{x_0}^{x_1} y(x) \sqrt{1 + y'(x)^2} dx.$$

Here the curve begins at $y(x_0) = \mathbf{y}_0$ and ends at $y(x_1) = \mathbf{y}_1$. If it is possible to pass a catenary curve through these points, then it describes the surface with minimum area. The calculus of variations offers an easy route to this conclusion. ■

Example 17.2.3 *Passage of Light through an Inhomogeneous Medium*

If we look at an object close to the horizon but well above the earth's surface, the light from it will bend as it passes through the atmosphere. This is a consequence of the fact that the speed of light decreases as it passes through an increasingly dense medium. If we assume the earth is flat and the speed of light $v(y)$ varies with the distance y above the earth, then the ray will take the path $y(x)$ of least time. The total travel time is

$$T(\mathbf{y}) = \int_{x_0}^{x_1} \frac{\sqrt{1 + y'(x)^2}}{v[y(x)]} dx$$

when the source is situated at (x_0, y_0) and we are situated at (x_1, y_1) . The calculus of variations provides theoretical insight into this generalization of Snell's problem. ■

Example 17.2.4 *Lagrange's versus Newton's Version of Mechanics*

The calculus of variations offers an alternative approach to classical mechanics. For a particle with kinetic energy $T(\mathbf{v}) = \frac{1}{2}m\|\mathbf{v}\|^2$ in a conservative force field with potential $U(x)$, we define the action integral

$$A(\mathbf{x}) = \int_{t_0}^{t_1} \{T[x'(t)] - U[x(t)]\} dt \quad (17.3)$$

on the path $x(t)$ of the particle from time t_0 to time t_1 . The path actually taken furnishes a stationary value of $A(\mathbf{x})$. One can demonstrate this fact by showing that the Euler-Lagrange equations coincide with Newton's equations of motion. ■

Example 17.2.5 *Isoperimetric Problem*

This classical Greek problem involves finding the plane curve of given length ℓ enclosing the greatest area. The circle of perimeter ℓ is the obvious solution. If we let a horizontal line segment form one side of the figure, then the solution is a circular arc. This version of the problem can be formalized by writing the enclosed area as

$$A(\mathbf{y}) = \int_{x_0}^{x_1} y(x) dx$$

and its length as

$$L(\mathbf{y}) = \int_{x_0}^{x_1} \sqrt{1 + y'(x)^2} dx.$$

The constrained problem of minimizing $A(\mathbf{y})$ subject to $L(\mathbf{y}) = \ell$ and $y(x_0) = y(x_1) = 0$ can be solved by introducing a Lagrange multiplier. ■

Example 17.2.6 *Splines*

A spline is a smooth curve interpolating a given function at specified points. From the variational perspective, we would like to find the function $x(t)$ minimizing the curvature

$$C(\mathbf{x}) = \int_a^b x''(t)^2 dt \tag{17.4}$$

subject to the constraints $x(s_i) = x_i$ at $n + 1$ points in the interval $[a, b]$. In this situation the solution has limited smoothness. The interpolation points are called nodes and are numbered so that $s_0 = a < s_1 < \dots < s_n = b$. ■

17.3 Normed Vector Spaces

In the calculus of variations, functions are viewed as vectors belonging to normed vector spaces of infinite dimension. The vector space and attached norm vary from problem to problem. Many of the concepts from finite-dimensional linear algebra extend without comment to the infinite-dimensional setting. The new concepts that crop up are fairly subtle, so it is worth spending some time on concrete examples.

For instance, consider the collection $\mathcal{C}^0[a, b]$ of continuous real-valued functions defined on a compact interval $[a, b]$. It is clear that $\mathcal{C}^0[a, b]$ is closed

under addition and scalar multiplication and that the function $x(t) \equiv 0$ serves as the zero vector or origin. The choice of norm is less obvious. Three natural possibilities are:

$$\begin{aligned}\|\mathbf{x}\|_\infty &= \sup_{t \in [a,b]} |x(t)| \\ \|\mathbf{x}\|_1 &= \int_a^b |x(t)| dt \\ \|\mathbf{x}\|_2 &= \left[\int_a^b |x(t)|^2 dt \right]^{1/2}.\end{aligned}$$

It is straightforward to check that $\|\mathbf{x}\|_\infty$ and $\|\mathbf{x}\|_1$ satisfy the requirements of a norm as set down in Chap. 2. To prove that $\|\mathbf{x}\|_2$ qualifies as a norm, it is best to define it in terms of the inner product

$$\langle \mathbf{x}, \mathbf{y} \rangle = \int_a^b x(t)y(t)dt$$

as

$$\|\mathbf{x}\|_2 = \sqrt{\langle \mathbf{x}, \mathbf{x} \rangle}$$

and invoke the Cauchy-Schwarz inequality.

In contrast to finite-dimensional spaces, not all norms are equivalent on infinite-dimensional spaces. Consider the sequence of continuous functions

$$x_n(t) = \begin{cases} 1 - nt & t \in [0, 1/n] \\ 0 & t \notin [0, 1/n] \end{cases}$$

on the unit interval. It is clear that $\|\mathbf{x}_n\|_\infty = 1$ for all n while

$$\begin{aligned}\|\mathbf{x}_n\|_1 &= \int_0^{1/n} (1 - nt) dt = \frac{1}{2n} \\ \|\mathbf{x}_n\|_2^2 &= \int_0^{1/n} (1 - nt)^2 dt = \frac{1}{3n}.\end{aligned}$$

Thus, $x_n(t)$ converges to the origin in two of these norms but not in the third.

Besides continuous functions, it is often useful to consider continuously differentiable functions. The vector space of such functions over the interval $[a, b]$ is denoted by $\mathcal{C}^1[a, b]$. The norm

$$\|\mathbf{x}\|_{\infty 1} = \sup_{t \in [a,b]} |x(t)| + \sup_{t \in [a,b]} |x'(t)| \quad (17.5)$$

is fairly natural when one is interested in uniform convergence of a sequence $x_n(t)$ together with its derivatives $x'_n(t)$. On the vector space $\mathcal{C}^m[a, b]$ of functions with m continuous derivatives, one can define the similar norm

$$\|\mathbf{x}\|_{\infty m} = \sum_{k=0}^m \sup_{t \in [a, b]} |x^{(k)}(t)|. \quad (17.6)$$

In most applications, the completeness of a normed vector space is an issue. A sequence $x_n(t)$ is said to be Cauchy if for every $\epsilon > 0$ there exists an integer n such that $\|\mathbf{x}_j - \mathbf{x}_k\| < \epsilon$ whenever $j \geq n$ and $k \geq n$. A normed vector space is complete if every Cauchy sequence possesses a limit in the space. For example, the vector space $\mathcal{C}^0[a, b]$ is complete under the uniform norm $\|\mathbf{x}\|_{\infty}$. For a proof of this fact, observe that

$$|x_j(t) - x_k(t)| \leq \|\mathbf{x}_j - \mathbf{x}_k\|_{\infty}$$

for every $t \in [a, b]$. This implies that the sequence $x_n(t)$ is Cauchy on the real line and possesses a limit $x(t)$. Since the convergence to $x(t)$ is uniform in t , Proposition 2.8.1 can be invoked to finish the proof.

The space $\mathcal{C}^0[a, b]$ is not complete under either of the norms $\|\mathbf{x}\|_1$ or $\|\mathbf{x}\|_2$. It is possible to extend $\mathcal{C}^0[a, b]$ to larger normed vector spaces $L_1[a, b]$ and $L_2[a, b]$ that are complete under these norms. The process of completion is one of the most fascinating parts of the theory of integration. Unfortunately, the work involved is far more than we can undertake here. Complete normed vector spaces are called Banach spaces; complete inner product spaces are called Hilbert spaces.

The space $\mathcal{C}^1[a, b]$ under the norm (17.5) is also complete. A little thought makes it clear that a Cauchy sequence $x_n(t)$ in $\mathcal{C}^1[a, b]$ not only converges uniformly to a continuous function $x(t)$, but its sequence of derivatives $x'_n(t)$ also converges uniformly to a continuous function $y(t)$. Applying the dominated convergence theorem to the sequence

$$x_n(t) = \int_a^t x'_n(s) ds$$

shows that

$$x(t) = \int_a^t y(s) ds.$$

It now follows from the fundamental theorem of calculus that $x'(t)$ exists and equals $y(t)$. A slight extension of this argument demonstrates that $\mathcal{C}^m[a, b]$ is complete under the norm (17.6). For this reason, we will tacitly assume in the sequel that $\mathcal{C}^m[a, b]$ is equipped with the $\|\mathbf{x}\|_{\infty m}$ norm.

17.4 Linear Operators and Functionals

Linear algebra focuses on linear maps and their matrix representations. In infinite-dimensional spaces, linear maps are referred to as linear operators. When the range of a linear operator is the real line, the operator is said to be a linear functional. Unfortunately, linear operators on normed vector spaces are no longer automatically continuous. Continuity is intimately tied to boundedness. A linear operator A from a normed linear space X with norm $\|\mathbf{x}\|_p$ to a normed linear space Y with norm $\|\mathbf{y}\|_q$ is bounded if there exists a constant c such that $\|A(\mathbf{x})\|_q \leq c\|\mathbf{x}\|_p$ for all $\mathbf{x} \in X$. The least such constant c determines the induced operator norm $\|A\|$. This verbal description just recapitulates the definition given in equation (2.2) of Chap. 2. For linear functionals, mathematicians invariably use the norm $\|y\|_q = |y|$ derived from the absolute value function. It is trivial to check that the collection of bounded linear operators between two normed vector spaces is closed under pointwise addition and scalar multiplication. Thus, this collection is a normed vector space in its own right.

Here are three typical bounded linear functionals:

$$A_1(\mathbf{x}) = x(d), \quad A_2(\mathbf{x}) = \int_a^b x(t)dt, \quad A_3(\mathbf{x}) = x'(d). \quad (17.7)$$

The first two of these are defined on the space $\mathcal{C}^0[a, b]$ and the third on the space $\mathcal{C}^1[a, b]$. The evaluation point d can be any point from $[a, b]$. Straightforward arguments show that the induced norms satisfy the inequalities $\|A_1\| \leq 1$, $\|A_2\| \leq b - a$, and $\|A_3\| \leq 1$.

Bounded linear operators are a little more exotic. If $y(t)$ is a monotone function mapping $[a, b]$ into itself, then the linear operator

$$A_4(\mathbf{x}) = \mathbf{x} \circ \mathbf{y}$$

composing \mathbf{x} and \mathbf{y} maps $\mathcal{C}^0[a, b]$ into itself. For example, if $[a, b] = [0, 1]$ and $y(t) = t^2$, then $A_4(\mathbf{x})(t) = x(t^2)$. The operator A_4 has norm $\|A_4\| \leq 1$. Because integration turns one continuous function into another, the linear operator

$$A_5(\mathbf{x})(s) = \int_a^s x(t)dt$$

also maps $\mathcal{C}^0[a, b]$ into itself. This operator has norm $\|A_5\| \leq b - a$. It is a special case of the linear operator

$$A_6(\mathbf{x})(s) = \int_a^b K(s, t)x(t)dt \quad (17.8)$$

defined by a bounded function $K(s, t)$ with domain the square $[a, b] \times [a, b]$. If $|K(s, t)| \leq c$ for all s and t , then the inequality

$$\begin{aligned} \left| \int_a^b K(s, t)x(t)dt \right| &\leq \int_a^b |K(s, t)||x(t)|dt \\ &\leq c\|x\|_\infty(b-a), \end{aligned}$$

shows that $\|A_6\| \leq c(b-a)$.

For a final example, let us define the injection operator $\text{Inj}_1(x)$ by the formula

$$\text{Inj}_1(\mathbf{x}) = \begin{pmatrix} \mathbf{x} \\ \mathbf{x}' \end{pmatrix}.$$

This is a linear operator from the normed vector space $\mathcal{C}^1[a, b]$ to the vector space $\mathcal{C}_2^0[a, b]$ of continuous vector-valued functions with two components. It is easy to check that

$$\left\| \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} \right\| = \|\mathbf{x}\|_\infty + \|\mathbf{y}\|_\infty$$

defines a norm on $\mathcal{C}_2^0[a, b]$. Furthermore, under this norm and the standard norm $\|\cdot\|_\infty$ on $\mathcal{C}^1[a, b]$, the linear operator $\text{Inj}_1(x)$ has induced operator norm 1. ■

The next proposition clarifies the relationship between boundedness and continuity.

Proposition 17.4.1 *The following three assertions concerning a linear operator A are equivalent:*

- (a) A is continuous,
- (b) A is continuous at the origin $\mathbf{0}$,
- (c) A is bounded.

Proof: Assertion (a) clearly implies assertion (b). Let the domain of A have norm $\|\cdot\|_p$ and the range norm $\|\cdot\|_q$. Assertion (c) implies assertion (a) because of the inequality

$$\|A(\mathbf{y}) - A(\mathbf{x})\|_q = \|A(\mathbf{y} - \mathbf{x})\|_q \leq \|A\| \cdot \|\mathbf{y} - \mathbf{x}\|_p.$$

To complete the proof, we must show that assertion (b) implies assertion (c). If A is unbounded, then there exists a sequence $\mathbf{x}_n \neq \mathbf{0}$ with

$$\|A(\mathbf{x}_n)\|_q \geq n\|\mathbf{x}_n\|_p.$$

If we set

$$\mathbf{y}_n = \frac{1}{n\|\mathbf{x}_n\|_p} \mathbf{x}_n,$$

then \mathbf{y}_n converges to $\mathbf{0}$, but $\|A(\mathbf{y}_n)\|_q \geq 1$ does not converge to 0. ■

17.5 Differentials

Our approach to differentiation is to replace slope matrices by slope operators. Let $F(\mathbf{y})$ be a nonlinear operator from a normed vector space U with norm $\|\cdot\|_p$ to a normed vector space V with norm $\|\cdot\|_q$. The right-hand side of the slope equation

$$F(\mathbf{y}) - F(\mathbf{x}) = S(\mathbf{y}, \mathbf{x})(\mathbf{y} - \mathbf{x}) \quad (17.9)$$

now involves a bounded linear operator $S(\mathbf{y}, \mathbf{x})$ from U to V operating on the vector $\mathbf{y} - \mathbf{x}$ in U . The operator $S(\mathbf{y}, \mathbf{x})$ has induced norm

$$\|S(\mathbf{y}, \mathbf{x})\| = \sup_{\|\mathbf{u}\|_p=1} \|S(\mathbf{y}, \mathbf{x})\mathbf{u}\|_q.$$

The operator $F(\mathbf{y})$ is said to have differential $dF(\mathbf{x})$ at \mathbf{x} provided the slope equation (17.9) holds for all \mathbf{y} sufficiently close to \mathbf{x} and

$$\lim_{\mathbf{y} \rightarrow \mathbf{x}} \|S(\mathbf{y}, \mathbf{x}) - dF(\mathbf{x})\| = 0.$$

This last equation implicitly requires $dF(\mathbf{x})$ to be a bounded linear operator from U to V . Furthermore, the affine map $\mathbf{y} \mapsto F(\mathbf{x}) + dF(\mathbf{x})(\mathbf{y} - \mathbf{x})$ uniformly approximates $F(\mathbf{y})$ in the sense that

$$\begin{aligned} \|F(\mathbf{y}) - F(\mathbf{x}) - dF(\mathbf{x})(\mathbf{y} - \mathbf{x})\|_q &= \|[S(\mathbf{y}, \mathbf{x}) - dF(\mathbf{x})](\mathbf{y} - \mathbf{x})\|_q \\ &\leq \|S(\mathbf{y}, \mathbf{x}) - dF(\mathbf{x})\| \cdot \|\mathbf{y} - \mathbf{x}\|_p \end{aligned}$$

for \mathbf{y} close to \mathbf{x} . Thus, we arrive at an appropriate extension of Carathéodory's definition of the differential. It is also possible to define Fréchet's differential in this setting. As Proposition 4.4.1 of Chap. 4 shows, the two definitions are equivalent on inner product spaces. On more general normed vector spaces, it is unclear whether a Fréchet differentiable function is necessarily Carathéodory differentiable.

The various rules for combining differentials remain in force, and versions of the inverse and implicit function theorems continue to hold. The proofs of these theorems must be modified to avoid an appeal to compactness [46]. Rather than concentrate on extending our earlier proofs, we prefer to offer some concrete examples. The simplest example is a bounded linear operator $F(\mathbf{x})$. In this case, application of the definition shows that $dF(\mathbf{x})\mathbf{u} = F(\mathbf{u})$. Here are some more subtle examples.

Example 17.5.1 The Squaring Operator

Consider the operator $F(\mathbf{y}) = \mathbf{y}^2$ on $C^0[a, b]$. The relation

$$F(\mathbf{y}) - F(\mathbf{x}) = (\mathbf{y} + \mathbf{x})(\mathbf{y} - \mathbf{x})$$

suggests that we take $S(\mathbf{y}, \mathbf{x}) = \mathbf{y} + \mathbf{x}$. For this to make sense, we reinterpret each of the symbols \mathbf{y} and \mathbf{x} as multiplication by the function in question. Thus, the operator \mathbf{y} takes $\mathbf{u} \in \mathcal{C}^0[a, b]$ to the function $y(t)u(t)$. The obvious bound $\|\mathbf{y}\mathbf{u}\|_\infty \leq \|\mathbf{y}\|_\infty \|\mathbf{u}\|_\infty$ says \mathbf{y} viewed as a linear operator has induced norm $\|\mathbf{y}\| \leq \|\mathbf{y}\|_\infty$. Once we apply \mathbf{y} to itself, then it becomes clear that $\|\mathbf{y}\| = \|\mathbf{y}\|_\infty$. Given these preliminaries, the identities

$$\|S(\mathbf{y}, \mathbf{x}) - 2\mathbf{x}\| = \|\mathbf{y} - \mathbf{x}\| = \|\mathbf{y} - \mathbf{x}\|_\infty$$

now demonstrate that $dF(\mathbf{x})$ exists and equals the multiplication operator $2\mathbf{x}$. ■

Example 17.5.2 *An Integration Functional*

Suppose the continuous function $f(t, x)$ has a continuous partial derivative $\partial_2 f(t, x) = \frac{\partial}{\partial x} f(t, x)$. If we fix t and choose the canonical slope function

$$s(t, y, x) = \int_0^1 \partial_2 f[t, x + s(y - x)] ds$$

of $f(t, x)$ as a function of x , then

$$f(t, y) - f(t, x) = s(t, y, x)(y - x).$$

Furthermore, this choice ensures that $s(t, y, x)$ is jointly continuous in its three arguments. Now consider the functional

$$F(\mathbf{x}) = \int_a^b f[t, x(t)] dt$$

on $\mathcal{C}^0[a, b]$. The equation

$$F(\mathbf{y}) - F(\mathbf{x}) = \int_a^b s[t, y(t), x(t)][y(t) - x(t)] dt$$

suggests the candidate differential

$$dF(\mathbf{x})\mathbf{u} = \int_a^b \partial_2 f[t, x(t)]u(t) dt. \quad (17.10)$$

To prove this contention, we first show that the linear functional

$$S(\mathbf{y}, \mathbf{x})\mathbf{u} = \int_a^b s[t, y(t), x(t)]u(t) dt$$

is bounded. Because the interval $[a, b]$ is compact, $x(t)$ is bounded. If we fix $\delta > 0$ and limit attention to those \mathbf{y} with $\|\mathbf{y} - \mathbf{x}\|_\infty < \delta$, then all values of $x(t)$ and $y(t)$ occur within an interval $[c, d]$. For such \mathbf{y} we have

$$\begin{aligned} \left| \int_a^b s[t, y(t), x(t)]u(t) dt \right| &\leq \int_a^b |s[t, y(t), x(t)]| \cdot |u(t)| dt \\ &\leq k(b - a)\|\mathbf{u}\|_\infty, \end{aligned}$$

where k is the supremum of $|s(t, y, x)|$ on $J = [a, b] \times [c, d] \times [c, d]$. This settles the question of boundedness. Given the uniform continuity of $s(t, y, x)$ on J , if we choose $\|\mathbf{y} - \mathbf{x}\|_\infty$ small enough, then $|s[t, y(t), x(t)] - s[t, x(t), x(t)]|$ can be made uniformly smaller than a preassigned $\epsilon > 0$. For those \mathbf{y} it follows that

$$\left| \int_a^b \{s[t, y(t), x(t)] - s[t, x(t), x(t)]\}u(t) dt \right| \leq \epsilon(b - a)\|\mathbf{u}\|_\infty,$$

and this implies that $S(\mathbf{y}, \mathbf{x})$ converges to $S(\mathbf{x}, \mathbf{x})$ in the relevant operator norm.

In the classical theory, equation (17.10) is viewed as a directional derivation. Fixing the “direction” $u(t)$ in function space, one calculates the directional derivative

$$\lim_{\epsilon \rightarrow 0} \frac{F(\mathbf{x} + \epsilon \mathbf{u}) - F(\mathbf{x})}{\epsilon} = \int_a^b \partial_2 f[t, x(t)]u(t) dt$$

by differentiation under the integral sign and the chain rule. This is rigorous as far as it goes, but it does not prove the existence of the differential. Having the full apparatus of differentials at our disposal unifies the theory and eases the process of generalization. ■

Example 17.5.3 *Differentials of More General Functionals*

Many of the classical examples of the calculus of variations are slight elaborations of the last example. For instance, we can replace the argument $x(t)$ of $F(\mathbf{x})$ by a continuous vector-valued function on $[a, b]$. The differential (17.10) is still valid provided we interpret $\partial_2 f(t, \mathbf{x})$ as the differential of $f(t, \mathbf{x})$ with respect to \mathbf{x} and assume $x(t)$ belongs to the normed vector space $\mathcal{C}_m^0[a, b]$ of continuous functions with m components for some m . Another profitable extension is to consider functionals of the form

$$G(\mathbf{x}) = \int_a^b f[t, x(t), x'(t)] dt$$

depending on $x'(t)$ as well as $x(t)$. Straightforward extension of our previous arguments yield

$$dG(\mathbf{x})\mathbf{u} = \int_a^b \{\partial_2 f[t, x(t), x'(t)]u(t) + \partial_3 f[t, x(t), x'(t)]u'(t)\} dt, \quad (17.11)$$

where $\partial_3 f(t, x, x') = \frac{\partial}{\partial x'} f(t, x, x')$. Now we must assume that both partial derivatives $\partial_2 f(t, x, x')$ and $\partial_3 f(t, x, x')$ are jointly continuous in all variables. Alternatively, we can derive this result by noting that $G(\mathbf{x})$ is the composition of the functional $F(\mathbf{x})$ of the last example with the injection operator $\text{Inj}_1(\mathbf{x})$. The differential (17.11) then reduces to the chain rule

$$dG(\mathbf{x}) = dF[\text{Inj}_1(\mathbf{x})]d\text{Inj}_1(\mathbf{x}).$$

This argument remains valid when $x(t)$ is vector-valued provided we interpret the partial derivatives $\partial_2 f(t, \mathbf{x}, \mathbf{x}')$ and $\partial_3 f(t, \mathbf{x}, \mathbf{x}')$ as partial differentials. Even this formula can be generalized by considering functionals

$$G(\mathbf{x}) = \int_a^b f[t, x(t), x'(t), \dots, x^{(k)}(t)] dt$$

depending on $x(t)$ and its first k derivatives. In this case, the formula

$$dG(\mathbf{x})\mathbf{u} = \int_a^b \sum_{j=0}^k \partial_{j+1} f[t, x(t), \dots, x^{(k)}(t)] u^{(j)}(t) dt \quad (17.12)$$

just summarizes the chain rule

$$dG(\mathbf{x}) = dF[\text{Inj}_k(\mathbf{x})]d\text{Inj}_k(\mathbf{x})$$

involving the injection operator $\text{Inj}_k(\mathbf{x}) = (x, x', \dots, x^{(k)})^*$. ■

17.6 The Euler-Lagrange Equation

In proving the Euler-Lagrange equation (17.2), we assume that the function $x(t)$ minimizes the functional (17.1) among all competing functions in the space $\mathcal{C}^1[a, b]$. When the boundary values $x(a) = c$ and $x(b) = d$ are fixed, we consider the revised functional

$$F(\mathbf{x} + \mathbf{u}) = \int_a^b f[t, x(t) + u(t), x'(t) + u'(t)] dt \quad (17.13)$$

defined on the set of continuously differentiable functions $u(t)$ with $u(a) = 0$ and $u(b) = 0$. This closed subspace $\mathcal{D}^1[a, b]$ of $\mathcal{C}^1[a, b]$ qualifies as a Banach space in its own right, and if $x(t)$ minimizes the original functional, then $\mathbf{u} = \mathbf{0}$ minimizes the revised functional (17.13).

Fermat's principle, which is valid on any normed vector space, requires the differential of $F(\mathbf{x} + \mathbf{u})$ to vanish at $\mathbf{u} = \mathbf{0}$. In view of Example (17.5.3), this means

$$dF(\mathbf{x})\mathbf{u} = \int_a^b \{\partial_2 f[t, x(t), x'(t)]u(t) + \partial_3 f[t, x(t), x'(t)]u'(t)\} dt, \quad (17.14)$$

must vanish for every $\mathbf{u} \in \mathcal{D}^1[a, b]$. The appearance of $u'(t)$ in addition to $u(t)$ in this last equation is awkward. We can eliminate it invoking the integration by parts formula

$$\int_a^b \partial_3 f[t, x(t), x'(t)]u'(t) dt = - \int_a^b \frac{d}{dt} \partial_3 f[t, x(t), x'(t)]u(t) dt \quad (17.15)$$

using the boundary conditions $u(a) = u(b) = 0$. If we put

$$v(t) = \partial_2 f[t, x(t), x'(t)] - \frac{d}{dt} \partial_3 f[t, x(t), x'(t)],$$

then Fermat's principle reads

$$\int_a^b v(t)u(t) dt = 0$$

for every $u \in \mathcal{D}^1[a, b]$. The special case $k = 0$ of the next lemma completes the proof of the Euler-Lagrange equation (17.2).

Proposition 17.6.1 (Du Bois-Reymond) *Let the function $v(t)$ be continuous on $[a, b]$ and satisfy*

$$\int_a^b v(t)u(t) dt = 0$$

for every $u(t)$ in $\mathcal{C}^k[a, b]$ with

$$u^{(j)}(a) = u^{(j)}(b) = 0, \quad 0 \leq j \leq k.$$

Then $v(t) = 0$ for all t .

Proof: Suppose $v(t)$ is not identically 0. By continuity there exists an interval $[c, d] \subset [a, b]$ on which $v(t)$ is either strictly positive or strictly negative. Without loss of generality, we assume the former case and take $a < c < d < b$. It is straightforward to prove by induction that the function

$$u(t) = \begin{cases} (t-c)^{k+1}(d-t)^{k+1} & t \in [c, d] \\ 0 & t \notin [c, d] \end{cases}$$

is continuously differentiable up to order k . Because the continuous function $v(t)u(t)$ is positive throughout the open interval (c, d) and vanishes outside it, the integral

$$\int_a^b v(t)u(t) dt > 0.$$

This contradiction proves the claim. ■

Two cases of the Euler-Lagrange equation (17.2) merit special mention. First, if $f(t, x, x')$ does not depend on x , then

$$\frac{d}{dt} \frac{\partial}{\partial x'} f[t, x(t), x'(t)] = 0,$$

and

$$\frac{\partial}{\partial x'} f[t, x(t), x'(t)] = c \tag{17.16}$$

for some constant c . Second, if $f(t, x, x')$ does not depend on t , then

$$x'(t) \frac{\partial}{\partial x'} f[t, x(t), x'(t)] - f[t, x(t), x'(t)] = c \quad (17.17)$$

is constant. This assertion follows from the identity

$$\begin{aligned} \frac{d}{dt} [x' \partial_3 f - f] &= x'' \partial_3 f + x' \frac{d}{dt} \partial_3 f - \partial_1 f - \partial_2 f x' - \partial_3 f x'' \\ &= x' \left[\frac{d}{dt} \partial_3 f - \partial_2 f \right] \\ &= 0. \end{aligned}$$

In the absence of fixed boundary values $x(a) = c$ and $x(b) = d$, the integration by parts argument invoked in deriving the Euler-Lagrange equations is still valid. However, perturbations with $u(a) \neq 0$ or $u(b) \neq 0$ are now pertinent. To ensure that the boundary contributions

$$\partial_3 f[b, x(b), x'(b)]u(b) - \partial_3 f[a, x(a), x'(a)]u(a)$$

to $dF(\mathbf{x})\mathbf{u}$ vanish, we must assume that the multiplier $\partial_3 f[a, x(a), x'(a)]$ of $u(a)$ vanishes when $x(a)$ is not fixed and the multiplier $\partial_3 f[b, x(b), x'(b)]$ of $u(b)$ vanishes when $x(b)$ is not fixed. These constraints are referred to as free or natural boundary conditions.

Some applications involve optimization over multivariate functions $x(t)$. In this case the Euler-Lagrange equation (17.2) holds component by component. Derivation of this result relies on considering multivariate perturbations $u(t)$ with all but one component identically 0. Our discussion of Lagrangian mechanics in the next section requires multivariate functions. Functionals depending on derivatives beyond the first derivative can also be treated by the methods described in this section as noted in Problem 15. To derive the appropriate generalization of the Euler-Lagrange equations, we again use integration by parts and Proposition 17.6.1. Our consideration of splines provides insight into how one deals with functionals depending on second derivatives.

The classical theory of the calculus of variations involves only necessary conditions for an optimum. The modern theory takes up sufficient conditions as well [102, 139]. Although it is impossible to do justice to the modern theory in a brief exposition, it is fair to point out the important role of convexity. A functional $F(\mathbf{x})$ is convex if Jensen's inequality

$$F[\alpha \mathbf{x} + (1 - \alpha)\mathbf{y}] \leq \alpha F(\mathbf{x}) + (1 - \alpha)F(\mathbf{y})$$

holds for all \mathbf{x} , \mathbf{y} , and $\alpha \in [0, 1]$. Proposition 6.5.2 forges the most helpful connection between convexity and optimization. When \mathbf{x} is a stationary point of $F(\mathbf{x})$, then the operative inequality $dF(\mathbf{x})\mathbf{u} \geq 0$ of the proposition automatically holds for all admissible u . A more crucial point is to

recognize convexity when it occurs. Fortunately, one simple test works. If the integrand $f(t, x, x')$ in the functional (17.1) is convex in its last two variables x and x' for each fixed t , then it is trivial to show that $F(\mathbf{x})$ is convex in \mathbf{x} . As pointed out momentarily, this is the case with geodesic problems.

17.7 Applications of the Euler-Lagrange Equation

To illustrate the solution of the Euler-Lagrange equations, we revisit the first four examples of Sect. 17.2. None of these examples involves constraints beyond fixed boundary conditions.

Example 17.7.1 Geodesic on a Cylinder

Without loss of generality, we suppose that the z axis and the central axis of the cylinder coincide. If the cylinder has radius r , then a curve on its surface can be represented by the triple $[r \cos \theta, r \sin \theta, z(\theta)]$ using cylindrical coordinates (θ, z) . The distance connecting the point (θ_0, z_0) to the point (θ_1, z_1) is

$$G(z) = \int_{\theta_0}^{\theta_1} \sqrt{r^2 + z'(\theta)^2} d\theta.$$

Here we assume $\theta_0 < \theta_1$ and $z_0 \neq z_1$. If $\theta_0 = \theta_1$, then it is clear that the geodesic is a vertical line on the surface. If $z_0 = z_1$, then the geodesic is a circular arc in a plane perpendicular to the central axis.

Because the function $f(\theta, z, z') = \sqrt{r^2 + z'^2}$ does not depend on z , the simplified Euler-Lagrange equation (17.16) applies. This requires

$$\frac{z'}{\sqrt{r^2 + z'^2}}$$

to be constant, which is achieved by taking z' to be constant. In view of the constraints $z(\theta_0) = z_0$ and $z(\theta_1) = z_1$, we have

$$z(\theta) = \frac{z_1 - z_0}{\theta_1 - \theta_0} \theta + \frac{\theta_1 z_0 - \theta_0 z_1}{\theta_1 - \theta_0}.$$

To prove that this solution provides the minimum of $G(z)$, it suffices to note that $f(\theta, z, z') = \sqrt{r^2 + z'^2}$ is convex in (z, z') . Convexity follows from the relationship between $\sqrt{r^2 + z'^2}$ and the Euclidean norm of \mathbb{R}^2 . ■

Example 17.7.2 Minimal Surface Area of Revolution

Because $f(x, y, y') = 2\pi y \sqrt{1 + y'^2}$ does not depend on x , the simplified Euler-Lagrange equation (17.17) requires

$$\frac{y'^2 y}{\sqrt{1 + y'^2}} - y \sqrt{1 + y'^2} = c$$

to be constant. Straightforward algebra now gives $y^2 = c^2(1 + y'^2)$ and

$$y' = \left[\left(\frac{y}{c} \right)^2 - 1 \right]^{1/2} = \frac{1}{c} \sqrt{y^2 - c^2}.$$

If we rewrite this as

$$\frac{dy}{\sqrt{y^2 - c^2}} = \frac{dx}{c},$$

then Table 4.1 shows that

$$\operatorname{arccosh}^{-1} \frac{y}{c} = \frac{x}{c} + d$$

for some constant d . It follows that

$$y(x) = c \cosh \left(\frac{x}{c} + d \right).$$

If we take $(x_0, y_0) = (0, a)$, then we can eliminate the constant c and express

$$y(x) = \frac{a}{\cosh d} \cosh \left(\frac{\cosh d}{a} x + d \right).$$

In some cases the adjustable parameter d can be used to match the other end of the catenary curve $y(x)$ to (x_1, y_1) . In other cases, there is no minimal surface of revolution in the ordinary sense [15]. ■

Example 17.7.3 *Passage of Light through an Inhomogeneous Medium*

This is another example of a integrand $f(x, y, y') = \sqrt{1 + y'^2}/v(y)$ that does not depend on x . The quantity (17.17)

$$\frac{y'^2}{v(y)\sqrt{1 + y'^2}} - \frac{\sqrt{1 + y'^2}}{v(y)} = -\frac{1}{v(y)\sqrt{1 + y'^2}}$$

is constant along the optimal path. If we differentiate the identity

$$-\ln v[y(x)] - \ln \sqrt{1 + y'^2(x)} = \ln(-c)$$

with respect to x , then the formula

$$y''(x) = -[1 + y'(x)^2] \frac{v'[y(x)]}{v[y(x)]}$$

emerges. Because the speed of light increases with increasing altitude, $v'[y(x)] > 0$, and hence $y''(x) < 0$. In other words, the path that light follows toward an observer on the earth is concave and bends downward. Hence, the setting sun appears higher in the sky than it actually is [240]. Problem 20 considers the special case $v(y) = y$, where one can explicitly solve for $y(x)$. ■

Example 17.7.4 *Lagrange's versus Newton's Version of Mechanics*

Let $x(t)$ be the path taken in \mathbb{R}^3 by a single particle with initial and final values $x(t_0) = \mathbf{a}$ and $x(t_1) = \mathbf{b}$. The Euler-Lagrange equations at the stationary value of the action integral (17.3) are

$$\begin{aligned} 0 &= -\frac{\partial}{\partial x_i}U[x(t)] - \frac{d}{dt}\frac{\partial}{\partial x'_i}T[x'(t)] \\ &= -\frac{\partial}{\partial x_i}U[x(t)] - mx''_i(t). \end{aligned}$$

These are clearly Newton's equations of motions in a conservative force field. The total energy $E = T(\mathbf{x}') + U(\mathbf{x})$ of the particle is conserved because

$$\frac{d}{dt}E(t) = m \sum_{i=1}^3 x'_i(t)x''_i(t) + \sum_{i=1}^3 \frac{\partial}{\partial x_i}U[x(t)]x'_i(t) = 0.$$

In many applications we follow several particles. For instance in the n -body problem, n particles move under the influence of their mutual gravitational fields. Let the i th particle have mass m_i and position $x_i(t)$ in \mathbb{R}^3 at time t . Assuming the gravitational constant is 1, the potential energy of the system is given by

$$U(\mathbf{x}) = - \sum_{\{i,j\}} \frac{m_i m_j}{\|\mathbf{x}_i - \mathbf{x}_j\|},$$

where the sum ranges over all pairs $\{i, j\}$ of particles and $x(t)$ is the vector constructed by stacking the $x_i(t)$. The total kinetic energy of the particles is

$$T(\mathbf{x}') = \frac{1}{2} \sum_{i=1}^n m_i \|\mathbf{x}'_i\|^2.$$

According to Lagrangian mechanics, the motion of the system yields the stationary value of the action integral (17.3). The Euler-Lagrange equations are

$$\begin{aligned} 0 &= -\frac{\partial}{\partial x_{ik}}U[x(t)] - \frac{d}{dt}\frac{\partial}{\partial x'_{ik}}T[x'(t)] \\ &= -\sum_{j \neq i} \frac{m_i m_j}{\|x_i(t) - x_j(t)\|^2} \cdot \frac{x_{ik}(t) - x_{jk}(t)}{\|x_i(t) - x_j(t)\|} - m_i x''_{ik}(t). \end{aligned}$$

In other words, the force exerted on particle i by particle j is given by Newton's inverse square law. The total energy $E = T(\mathbf{x}') + U(\mathbf{x})$ of the system is conserved because

$$\begin{aligned} \frac{d}{dt}E(t) &= \sum_{i=1}^n m_i \sum_{k=1}^3 x'_{ik}(t)x''_{ik}(t) + \sum_{i=1}^n \sum_{k=1}^3 \frac{\partial}{\partial x_{ik}}U[x(t)]x'_{ik}(t) \\ &= 0. \end{aligned}$$

This result holds for any conservative dynamical system governed by the Euler-Lagrange equations with potential $U(\mathbf{x})$. Conservation of linear and angular momentum also hold in the n -body framework as documented in Problem 21. Dynamical systems with more general potential functions do not necessarily preserve linear and angular momentum. ■

17.8 Lagrange's Lacuna

Let us now expose and correct a deception foisted on the reader in deriving the Euler-Lagrange equations. At a crucial point in our derivation we invoked the integration-by-parts formula (17.15) without proving that the factor $\partial_3 f[t, x(t), x'(t)]$ is continuously differentiable. This gap in Lagrange's original treatment was noted by Du Bois-Reymond. His correction is based on the following variant of Proposition 17.6.1.

Proposition 17.8.1 *Let the function $v(t)$ be continuous on $[a, b]$ and satisfy*

$$\int_a^b v(t)u^{(k)}(t) dt = 0$$

for every test function $u(t)$ in $\mathcal{C}^k[a, b]$ with

$$u^{(j)}(a) = u^{(j)}(b) = 0, \quad 0 \leq j \leq k-1.$$

Then $v(t)$ is a polynomial of degree $k-1$ on $[a, b]$.

Proof: Integration by parts shows that

$$\int_a^b p(t)u^{(k)}(t) dt = 0$$

and therefore that

$$\int_a^b [v(t) - p(t)]u^{(k)}(t) dt = 0$$

for every test function $u(t)$ and polynomial $p(t)$ of degree $k-1$. If we can construct a test function $u(t)$ and a polynomial $p(t)$ of degree $k-1$ such that $u^{(k)}(t) = v(t) - p(t)$, then the identity

$$\int_a^b [v(t) - p(t)]^2 dt = 0$$

implies $v(t) = p(t)$. Construction of $u(t)$ involves some technicalities that we will avoid by focusing on the case $k=2$. The proof of the case $k=1$ is

similar but simpler. If we put $p(t) = c + d(t - a)$, then it seems sensible to define

$$\begin{aligned} u'(t) &= \int_a^t [v(s) - c - d(s - a)] ds \\ &= \int_a^t v(s) ds - c(t - a) - \frac{d}{2}(t - a)^2 \\ u(t) &= \int_a^t \int_a^s v(r) dr ds - c \int_a^t (s - a) ds - \frac{d}{2} \int_a^t (s - a)^2 ds \\ &= \int_a^t \int_a^s v(r) dr ds - \frac{c}{2}(t - a)^2 - \frac{d}{6}(t - a)^3. \end{aligned}$$

Because $u'(a) = u(a) = 0$ is clearly true, the only remaining issue is whether we can choose c and d so that $u'(b) = u(b) = 0$. However, this is possible since the matrix implementing the linear system

$$\begin{aligned} 0 &= \int_a^b v(s) ds - c(b - a) - \frac{d}{2}(b - a)^2 \\ 0 &= \int_a^b \int_a^s v(r) dr ds - \frac{c}{2}(b - a)^2 - \frac{d}{6}(b - a)^3 \end{aligned}$$

is invertible. Indeed, invertibility follows from the identity

$$\det \begin{pmatrix} (b - a) & (b - a)^2/2 \\ (b - a)^2/2 & (b - a)^3/6 \end{pmatrix} = -(b - a)^4/12.$$

Thus, c and d can be chosen so that $u(t)$ satisfies the requisite boundary conditions. ■

To apply Proposition 17.8.1 in the derivation of the Euler-Lagrange equation, we define the function

$$g(t) = \int_a^t \partial_2 f[s, x(s), x'(s)] ds.$$

The fundamental theorem of calculus implies that $g(t)$ is continuously differentiable. Hence, integration by parts gives the alternative

$$dF(\mathbf{x})\mathbf{u} = \int_a^b \{-g(t) + \partial_3 f[t, x(t), x'(t)]\} u'(t) dt$$

to equation (17.14). According to Proposition 17.8.1 with $k = 1$, the function

$$-g(t) + \partial_3 f[t, x(t), x'(t)] = c$$

for some constant c . It follows that $\partial_3 f[t, x(t), x'(t)] = g(t) + c$ is continuously differentiable as required.

17.9 Variational Problems with Constraints

The isoperimetric problem requires maximizing the area enclosed by a curve of fixed perimeter. More generally suppose we wish to maximize a functional $F(\mathbf{y})$ subject to the equality constraints $G_i(\mathbf{y}) = 0$ for $1 \leq i \leq p$. If \mathbf{x} furnishes a local maximum, then we can hope that the Lagrange multiplier rule

$$dF(\mathbf{x})\mathbf{u} + \sum_{i=1}^p \lambda_i dG_i(\mathbf{x})\mathbf{u} = 0 \quad (17.18)$$

will be valid for all admissible functions u . This turns out to be the case if the differentials $dG_1(\mathbf{x}), \dots, dG_p(\mathbf{x})$ are linearly independent. Linear independence fails whenever there exists a nontrivial vector \mathbf{c} with components c_1, \dots, c_p such that

$$\sum_{i=1}^p c_i dG_i(\mathbf{x})\mathbf{u} = 0$$

for all possible \mathbf{u} . Such a failure is impossible if there exists a finite sequence $\mathbf{u}_1, \dots, \mathbf{u}_p$ of admissible functions such that the square matrix $[dG_i(\mathbf{x})\mathbf{u}_j]$ is invertible. For the sake of simplicity, we will take linear independence to mean that for some choice of $\mathbf{u}_1, \dots, \mathbf{u}_p$ the matrix $[dG_i(\mathbf{x})\mathbf{u}_j]$ is invertible. With this stipulation in mind, we now prove the multiplier rule (17.18) under independence.

Our strategy will be to examine map

$$H(\boldsymbol{\alpha}) = \begin{pmatrix} F(\mathbf{x} + \alpha_0\mathbf{u} + \alpha_1\mathbf{u}_1 + \cdots + \alpha_p\mathbf{u}_p) \\ G_1(\mathbf{x} + \alpha_0\mathbf{u} + \alpha_1\mathbf{u}_1 + \cdots + \alpha_p\mathbf{u}_p) \\ \vdots \\ G_p(\mathbf{x} + \alpha_0\mathbf{u} + \alpha_1\mathbf{u}_1 + \cdots + \alpha_p\mathbf{u}_p) \end{pmatrix}$$

defined for the given functions $\mathbf{u}_1, \dots, \mathbf{u}_p$ and an arbitrary admissible function \mathbf{u} . Note that $H(\boldsymbol{\alpha})$ maps \mathbb{R}^{p+1} into itself. The differential of $H(\boldsymbol{\alpha})$ at $\mathbf{0}$ is the Jacobian matrix

$$dH(\mathbf{0}) = \begin{pmatrix} dF(\mathbf{x})\mathbf{u} & dF(\mathbf{x})\mathbf{u}_1 & \cdots & dF(\mathbf{x})\mathbf{u}_p \\ dG_1(\mathbf{x})\mathbf{u} & dG_1(\mathbf{x})\mathbf{u}_1 & \cdots & dG_1(\mathbf{x})\mathbf{u}_p \\ \vdots & \vdots & \vdots & \vdots \\ dG_p(\mathbf{x})\mathbf{u} & dG_p(\mathbf{x})\mathbf{u}_1 & \cdots & dG_p(\mathbf{x})\mathbf{u}_p \end{pmatrix}.$$

Assuming the objective functional $F(\mathbf{x})$ and the constraints $G_i(\mathbf{x})$ are continuously differentiable in a neighborhood of \mathbf{x} , the function $H(\boldsymbol{\alpha})$ is continuously differentiable in a neighborhood of $\mathbf{0}$. Hence, we can invoke the inverse function theorem, Proposition 4.6.1, provided $dH(\mathbf{0})$ is invertible.

Assume for the sake of argument that this is the case. Then $H(\boldsymbol{\alpha})$ maps a neighborhood of $\mathbf{0}$ onto a neighborhood of the image

$$[F(\mathbf{x}), G_1(\mathbf{x}), \dots, G_p(\mathbf{x})]^* = [F(\mathbf{x}), 0, \dots, 0]^*.$$

Taking the neighborhood of $\mathbf{0}$ to be arbitrarily small implies that there are functions $\mathbf{v} = \mathbf{x} + \alpha_0 \mathbf{u} + \alpha_1 \mathbf{u}_1 + \dots + \alpha_p \mathbf{u}_p$ arbitrarily close to \mathbf{x} satisfying $F(\mathbf{v}) > F(\mathbf{x})$ and $G_i(\mathbf{v}) = 0$ for all i . This contradicts the assumption that \mathbf{x} furnishes a local maximum subject to the constraints. It follows that the matrix $dH(\mathbf{0})$ must be singular.

The connection of this condition to the multiplier rule (17.18) becomes less obscure when we exploit the fact that $\det dH(\mathbf{0}) = 0$. Expanding this determinant on the first column of $dH(\mathbf{0})$ leads to the equation

$$0 = \mu_0 dF(\mathbf{x})\mathbf{u} + \sum_{i=1}^p \mu_i dG_i(\mathbf{x})\mathbf{u}, \quad (17.19)$$

where the cofactor $\mu_0 = \det[dG_i(\mathbf{x})\mathbf{u}_j]$ is nonzero by assumption. If we divide equation (17.19) by μ_0 , then we arrive at the multiplier rule (17.18) with $\lambda_i = \mu_i/\mu_0$.

Example 17.9.1 Isoperimetric Problem

If we assume that the perimeter constraint has a nontrivial differential, then the combination of the multiplier rule and the Euler-Lagrange equation (17.2) implies

$$1 - \frac{d}{dx} \left[\lambda \frac{y'(x)}{\sqrt{1 + y'(x)^2}} \right] = 0.$$

This forces λ to be nonzero, and integration on x produces

$$\frac{y'(x)}{\sqrt{1 + y'(x)^2}} = cx + d$$

for $c = 1/\lambda$ and d a constant of integration. This last equation can be solved for $y'(x)$ in the form

$$y'(x) = \frac{cx + d}{\sqrt{1 - (cx + d)^2}}$$

and integrated to give

$$y(x) = -\frac{1}{c} \sqrt{1 - (cx + d)^2} + e$$

for some additional constant e . If we take $x_0 = -1$ and $x_1 = 1$, then application of the boundary conditions $y(-1) = y(1) = 0$ yields

$$\begin{aligned} e^2 &= \frac{1}{c^2} [1 - (-c + d)^2] \\ e^2 &= \frac{1}{c^2} [1 - (c + d)^2]. \end{aligned}$$

It follows that $d = 0$ and $e = \sqrt{c^{-2} - 1}$. Finally, the constraint

$$\begin{aligned} \ell &= \int_{-1}^1 \sqrt{1 + y'(x)^2} dx \\ &= \int_{-1}^1 \frac{1}{\sqrt{1 - c^2 x^2}} dx \\ &= \frac{2}{c} \arcsin c \end{aligned}$$

determines the constant c . More precisely, the line $c \mapsto \ell c$ and the convex increasing curve $c \mapsto 2 \arcsin c$ intersect in a single point in $(0, 1]$ whenever $2 < \ell \leq \pi$. Because

$$[y(x) - e]^2 + x^2 = \frac{1}{c^2},$$

the function $y(x)$ traces out a circular arc. If $\ell = \pi$, then the arc is a semicircle with center at the origin $\mathbf{0}$.

This analysis is predicated on the assumption that there exists a continuously differentiable function $u(x)$ with

$$\begin{aligned} dL(\mathbf{y})\mathbf{u} &= \int_{-1}^1 \frac{y'(x)}{\sqrt{1 + y'(x)^2}} u'(x) dx \\ &= \int_{-1}^1 cxu'(x) dx \\ &\neq 0 \end{aligned}$$

and $u(-1) = u(1) = 0$. One obvious choice is $u(x) = (1 - x^2)^2$. ■

17.10 Natural Cubic Splines

Our treatment of splines involves functions from the Banach space $\mathcal{C}^2[a, b]$. If $x(t)$ minimizes the spline functional (17.4), then the differential

$$dC(\mathbf{x})\mathbf{u} = 2 \int_a^b x''(t)u''(t) dt$$

is a special case of equation (17.12) and vanishes for all admissible $u(t)$. Such functions $u(t)$ satisfy $u(s_i) = 0$ at every node s_i . If we focus on test functions $u(t)$ that vanish outside a node interval $[s_{i-1}, s_i]$, then we can invoke Proposition 17.8.1 and infer that $x''(t)$ is linear on $[s_{i-1}, s_i]$. This compels $x(t)$ to be a piecewise cubic polynomial throughout $[a, b]$.

On node interval $[s_{i-1}, s_i]$, integration by parts and the cubic nature of $x(t)$ yield

$$\int_{s_{i-1}}^{s_i} x''(t)u''(t) dt = x''(t)u'(t)\Big|_{s_{i-1}}^{s_i} - x'''(t)u(t)\Big|_{s_{i-1}}^{s_i}. \quad (17.20)$$

The constraint $u(s_{i-1}) = u(s_i) = 0$ forces the boundary terms involving $x'''(t)$ to vanish. When we add the contributions (17.20) to form the full integral $\int_a^b x''(t)u''(t) dt$, all of the boundary terms involving $x''(t)$ cancel except for $x''(s_n)u'(s_n) - x''(s_0)u'(s_0)$. If we impose the natural boundary conditions $x''(s_n) = x''(s_0) = 0$, then all terms vanish, and we recover the necessary condition $dC(\mathbf{x})\mathbf{u} = 2 \int_a^b x''(t)u''(t) dt = 0$ for optimality.

It is possible to demonstrate that there is precisely one piecewise cubic polynomial $x(t)$ from $\mathcal{C}^2[a, b]$ that interpolates the given values x_i at the nodes s_i and satisfies the natural boundary conditions $x''(s_0) = x''(s_n) = 0$. This exercise in linear algebra is carried out in the reference [166]. It is shown there that $x(t)$ minimizes $C(\mathbf{y})$ subject to the interpolation constraints. The fact that the minimum is unique is hardly surprising given the convexity of the functional $C(\mathbf{y})$ on $\mathcal{C}^2[a, b]$.

17.11 Problems

1. Prove that the space $\mathcal{C}^0[a, b]$ is not finite dimensional. (Hint: The functions $1, t, \dots, t^n$ are linearly independent.)
2. Prove that the set $P_n[0, 1]$ of polynomials of degree n or less on $[0, 1]$ is a closed subspace of $\mathcal{C}^0[0, 1]$.
3. Suppose X is an inner product space. Prove that the induced norm satisfies the parallelogram identity

$$\|\mathbf{u} + \mathbf{v}\|^2 + \|\mathbf{u} - \mathbf{v}\|^2 = 2(\|\mathbf{u}\|^2 + \|\mathbf{v}\|^2)$$

for all vectors \mathbf{u} and \mathbf{v} . Show that the norm of $\mathcal{C}^0[0, 1]$ is not induced by an inner product by producing \mathbf{u} and \mathbf{v} that fail the parallelogram identity.

4. Demonstrate that the closed unit ball $B = \{\mathbf{x} : \|\mathbf{x}\| \leq 1\}$ in $\mathcal{C}^0[0, 1]$ is not compact. (Hint: The sequence $1, t, t^2, t^3, \dots$ has no uniformly convergent subsequence.)

5. Demonstrate that a closed subset of a Banach space is complete.
6. Let X and Y be normed vector spaces with norms $\|\mathbf{x}\|_p$ and $\|\mathbf{y}\|_q$. The product space $X \times Y$ is a vector space if addition and scalar multiplication are defined coordinatewise. Show that the following are equivalent norms on $X \times Y$:

$$\begin{aligned}\|(\mathbf{x}, \mathbf{y})\|_\infty &= \max\{\|\mathbf{x}\|_p, \|\mathbf{y}\|_q\} \\ \|(\mathbf{x}, \mathbf{y})\|_1 &= \|\mathbf{x}\|_p + \|\mathbf{y}\|_q \\ \|(\mathbf{x}, \mathbf{y})\|_2 &= \sqrt{\|\mathbf{x}\|_p^2 + \|\mathbf{y}\|_q^2}.\end{aligned}$$

See Example 2.5.6 for the notion of equivalent norms.

7. Continuing Problem 6, show that $X \times Y$ is a Banach space under any of three proposed norms whenever X and Y are both Banach spaces.
8. Demonstrate that the three linear functionals A_1 , A_2 , and A_3 defined in (17.7) actually have the operator norms 1, $b - a$, and 1 on their respective normed linear spaces.
9. Suppose the function $K(s, t)$ in equation (17.8) is square-integrable over $[a, b] \times [a, b]$. Prove that the corresponding operator A_6 maps the Hilbert space $L_2[a, b]$ into itself in such a way that

$$\|A_6\|^2 \leq \int_a^b \int_a^b K(s, t)^2 ds dt.$$

10. Let $L(\mathbf{y})$ be a continuous linear functional on a normed vector space. Under what circumstances does $L(\mathbf{y})$ have a minimum or a maximum?
11. On an inner product space prove that the functional $\|\mathbf{x}\|^2$ is differentiable with differential $dF(\mathbf{x})\mathbf{y} = 2\langle \mathbf{x}, \mathbf{y} \rangle$. Deduce from this fact that the norm $G(\mathbf{x}) = \|\mathbf{x}\|$ is differentiable with differential $dG(\mathbf{x})\mathbf{y} = \langle \frac{\mathbf{x}}{\|\mathbf{x}\|}, \mathbf{y} \rangle$ whenever $\mathbf{x} \neq \mathbf{0}$.
12. Consider the linear functional $F(\mathbf{y}) = \langle \mathbf{v}, \mathbf{y} \rangle$ on an inner product space. Here the vector \mathbf{v} is fixed. Let $g(t)$ be a differentiable function on the real line, and define $G(\mathbf{y}) = g \circ F(\mathbf{y})$. Show that $G(\mathbf{y})$ has differential $dG(\mathbf{x})\mathbf{u} = g'(\langle \mathbf{v}, \mathbf{x} \rangle)\langle \mathbf{v}, \mathbf{u} \rangle$ at the vector \mathbf{x} .
13. Some variational problems have no solution. For instance, show that the following two functionals have no minimum subject to the given constraints:

- (a) The integral $\int_0^1 \sqrt{1 + y'(x)^2} dx$ subject to the four constraints $y(0) = y(1) = 0$ and $y'(0) = y'(1) = 1$.

- (b) Weierstrass's integral $\int_{-1}^1 x^2 y'(x)^2 dx$ subject to the constraints $y(-1) = -1$ and $y(1) = 1$.

In each case $y(x)$ can be any piecewise differentiable function on the given interval.

14. Consider the closed subspace of $C^1[0, 2\pi]$ defined by functions that are periodic. Wirtinger's inequality says

$$\int_0^{2\pi} f'(t)^2 dt \geq \int_0^{2\pi} f(t)^2 dt$$

for every such function $f(t)$ with equality if and only if

$$f(t) = a \cos t + b \sin t.$$

Prove this result by considering the functional

$$W(\mathbf{f}) = \int_0^{2\pi} [f'(t)^2 - f(t)^2] dt$$

For the purpose of this problem you may assume that the minimum of $W(\mathbf{f})$ is attained.

15. Consider a functional

$$F(\mathbf{x}) = \int_a^b f[t, x(t), x'(t), \dots, x^{(k)}(t)] dt$$

defined on $C^k[a, b]$. Derive the Euler-Lagrange equation

$$\frac{\partial}{\partial x} f + \sum_{j=1}^k (-1)^j \frac{d^j}{dt^j} \left[\frac{\partial}{\partial x^{(j)}} f \right] = 0$$

for the function $x(t)$ optimizing $F(\mathbf{x})$. What assumptions are necessary to justify this derivation?

16. Let $y(x)$ denote a continuously differentiable curve over the interval $[-1, 1]$ satisfying $y(-1) = y(1) = 0$. Find $y(x)$ that minimizes the length

$$L(\mathbf{y}) = \int_{-1}^1 \sqrt{1 + y'(x)^2} dx.$$

while enclosing a fixed area

$$\int_{-1}^1 y(x) dx = A \leq \frac{\pi}{2}.$$

17. If possible, find the minimum and maximum of the functional

$$F(\mathbf{y}) = \int_0^1 [y'(x)^2 + x^2] dx$$

over $\mathcal{C}^1[0, 1]$ subject to the constraints $y(0) = 0$, $y(1) = 0$, and $\int_0^1 y(x)^2 dx = 1$. Verify that your minimum solution yields the minimum value.

18. Find the minimum value of the functional

$$M(\mathbf{y}) = \int_0^1 xy(x) dx$$

over $\mathcal{C}^0[0, 1]$ subject to the constraint

$$V(\mathbf{y}) = \int_0^1 y(x)^2 dx = \frac{1}{12}.$$

Verify that your solution gives the minimum.

19. Find the minimum value of the functional

$$F(\mathbf{y}) = \int_0^\pi [y'(x)^2 + 2y(x) \sin x] dx$$

subject to the single constraint $y(0) = 0$. How does this differ from the solution when the constraint $y(\pi) = 0$ is added? Verify in each case that the minimum is achieved.

20. Suppose the speed of light in the upper-half
- xy
- plane is
- $v(y) = y$
- . Find the path of light connecting
- (x_0, y_0)
- to
- (x_1, y_1)
- . Show that the travel time is infinite if and only if either
- y_0
- or
- y_1
- equals 0.

21. In the
- n
- body problem, linear and angular momentum are defined by

$$L(t) = \sum_{i=1}^n m_i x'_i(t)$$

$$A(t) = \sum_{i=1}^n m_i x_i(t) \times x'_i(t).$$

Prove that these quantities are conserved. (Hints: The force exerted by i on j is equal in magnitude and opposite in direction to the force exerted by j on i . The cross product $\mathbf{v} \times \mathbf{v}$ of a vector with itself vanishes.)

22. Suppose the integrand $\mathcal{L}(t, \mathbf{x}, \mathbf{v})$ of the functional

$$F(\mathbf{x}) = \int_{t_0}^{t_1} \mathcal{L}[t, x(t), x'(t)] dt$$

is strictly convex in the variable \mathbf{v} with continuous second differential $\partial_3^2 \mathcal{L}(t, \mathbf{x}, \mathbf{v})$. In mechanics, $\mathcal{L}(t, \mathbf{x}, \mathbf{v})$ is called the Lagrangian, and its Fenchel conjugate

$$\mathcal{H}(t, \mathbf{x}, \mathbf{p}) = \sup_{\mathbf{v}} [\mathbf{p}^* \mathbf{v} - \mathcal{L}(t, \mathbf{x}, \mathbf{v})]$$

is called the Hamiltonian. According to Proposition 14.3.2, we recover $\mathcal{L}(t, \mathbf{x}, \mathbf{v})$ as

$$\mathcal{L}(t, \mathbf{x}, \mathbf{v}) = \sup_{\mathbf{p}} [\mathbf{v}^* \mathbf{p} - \mathcal{H}(t, \mathbf{x}, \mathbf{p})].$$

Show that \mathbf{v} and \mathbf{p} determine each other through the equations

$$\begin{aligned} \partial_3 \mathcal{L}(t, \mathbf{x}, \mathbf{v}) &= \mathbf{p}^* \\ \partial_3 \mathcal{H}(t, \mathbf{x}, \mathbf{p}) &= \mathbf{v}^* \end{aligned}$$

and that

$$\mathcal{H}(t, \mathbf{x}, \mathbf{p}) = \mathbf{p}^* \mathbf{v} - \mathcal{L}(t, \mathbf{x}, \mathbf{v}) \quad (17.21)$$

when these equations are satisfied. From the expression (17.21) deduce that

$$\begin{aligned} d\mathcal{H}(t, \mathbf{x}, \mathbf{p}) &= -d\mathcal{L}(t, \mathbf{x}, \mathbf{v}) + \mathbf{p}^* d\mathbf{v} + \mathbf{v}^* d\mathbf{p} \\ &= -\partial_1 \mathcal{L}(t, \mathbf{x}, \mathbf{v}) dt - \partial_2 \mathcal{L}(t, \mathbf{x}, \mathbf{v}) d\mathbf{x} + \partial_3 \mathcal{H}(t, \mathbf{x}, \mathbf{p}) d\mathbf{p}. \end{aligned}$$

Compare this to the differential

$$d\mathcal{H}(t, \mathbf{x}, \mathbf{p}) = \partial_1 \mathcal{H}(t, \mathbf{x}, \mathbf{p}) dt + \partial_2 \mathcal{H}(t, \mathbf{x}, \mathbf{p}) d\mathbf{x} + \partial_3 \mathcal{H}(t, \mathbf{x}, \mathbf{p}) d\mathbf{p}$$

and conclude that

$$\partial_2 \mathcal{L}(t, \mathbf{x}, \mathbf{v}) = -\partial_2 \mathcal{H}(t, \mathbf{x}, \mathbf{p}).$$

Use this result to prove that the Euler-Lagrange equation

$$\partial_2 \mathcal{L}[t, x(t), x'(t)] = \frac{d}{dt} \partial_3 \mathcal{L}[t, x(t), x'(t)]$$

is equivalent to the two Hamiltonian equations

$$\begin{aligned} \frac{d}{dt} x(t) &= \partial_3 \mathcal{H}[t, x(t), p(t)]^* \\ \frac{d}{dt} p(t) &= -\partial_2 \mathcal{H}[t, x(t), p(t)]^*. \end{aligned}$$

Thus, the first-order Hamiltonian differential equations can serve as a substitute for the second-order Euler-Lagrange differential equation in these circumstances.

23. Continuing Problem 22, show that the Hamiltonian is constant over time t whenever it does not depend explicitly on time, that is

$$\mathcal{H}(t, \mathbf{x}, \mathbf{p}) = \mathcal{H}(\mathbf{x}, \mathbf{p}).$$

24. Continuing Problem 22, suppose $\mathcal{L}(t, \mathbf{x}, \mathbf{v})$ can be expressed as the difference between the kinetic energy $T(\mathbf{v}) = \frac{1}{2} \sum_{i=1}^n m_i \|\mathbf{v}_i\|^2$ and the potential energy $U(\mathbf{x})$. Demonstrate that $\mathcal{H}(t, \mathbf{x}, \mathbf{p})$ is the sum $T(\mathbf{v}) + U(\mathbf{x})$ of the kinetic and potential energies. (Hint: See Problem 10 of Chap. 14.)

Appendix: Mathematical Notes

A.1 Univariate Normal Random Variables

A random variable X is said to be standard normal if it possesses the density function

$$\psi(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}}.$$

To find the characteristic function $\hat{\psi}(s) = E(e^{isX})$ of X , we derive and solve a differential equation. Differentiation under the integral sign and integration by parts together imply that

$$\begin{aligned} \frac{d}{ds} \hat{\psi}(s) &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{isx} i x e^{-\frac{x^2}{2}} dx \\ &= -\frac{i}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{isx} \frac{d}{dx} e^{-\frac{x^2}{2}} dx \\ &= -\frac{i}{\sqrt{2\pi}} e^{isx} e^{-\frac{x^2}{2}} \Big|_{-\infty}^{\infty} - \frac{s}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{isx} e^{-\frac{x^2}{2}} dx \\ &= -s \hat{\psi}(s). \end{aligned}$$

The unique solution to this differential equation with initial value $\hat{\psi}(0) = 1$ is $\hat{\psi}(s) = e^{-s^2/2}$. The differential equation also yields the moments

$$E(X) = \frac{1}{i} \frac{d}{ds} \hat{\psi}(0) = 0$$

and

$$\begin{aligned} E(X^2) &= \frac{1}{i^2} \frac{d^2}{ds^2} \hat{\psi}(0) \\ &= \frac{1}{i^2} \left[-\hat{\psi}(s) + s^2 \hat{\psi}(s) \right]_{s=0} \\ &= 1. \end{aligned}$$

An affine transformation $Y = \sigma X + \mu$ of X is normally distributed with density

$$\frac{1}{\sigma} \psi\left(\frac{y - \mu}{\sigma}\right) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-\mu)^2}{2\sigma^2}}.$$

Here we take $\sigma > 0$. The general identity $E\left[e^{is(\mu+\sigma X)}\right] = e^{is\mu} E\left[e^{i(\sigma s)X}\right]$ permits us to write the characteristic function of Y as

$$e^{is\mu} \hat{\psi}(\sigma s) = e^{is\mu - \frac{\sigma^2 s^2}{2}}.$$

The mean and variance of Y are μ and σ^2 .

One of the most useful properties of normally distributed random variables is that they are closed under the formation of independent linear combinations. Thus, if Y and Z are independent and normally distributed, then $aY + bZ$ is normally distributed for any choice of the constants a and b . To prove this result, it suffices to assume that Y and Z are standard normal. In view of the form of $\hat{\psi}(s)$, we then have

$$E\left[e^{is(aY+bZ)}\right] = E\left[e^{i(as)Y}\right] E\left[e^{i(bs)Z}\right] = \hat{\psi}\left(\sqrt{a^2 + b^2}s\right).$$

Thus, if we accept the fact that a distribution function is uniquely defined by its characteristic function, $aY + bZ$ is normally distributed with mean 0 and variance $a^2 + b^2$.

Doubtless the reader is also familiar with the central limit theorem. For the record, recall that if X_n is a sequence of independent identically distributed random variables with common mean μ and common variance σ^2 , then

$$\lim_{n \rightarrow \infty} \Pr \left[\frac{\sum_{j=1}^n (X_j - \mu)}{\sqrt{n\sigma^2}} \leq x \right] = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-\frac{u^2}{2}} du.$$

Of course, there is a certain inevitability to the limit being standard normal; namely, if the X_n are standard normal to begin with, then the standardized sum $n^{-1/2} \sum_{j=1}^n X_j$ is also standard normal.

A.2 Multivariate Normal Random Vectors

We now extend the univariate normal distribution to the multivariate normal distribution. Among the many possible definitions, we adopt the one most widely used in stochastic simulation. Our point of departure will be random vectors with independent standard normal components. If such a random vector \mathbf{X} has n components, then its density is

$$\prod_{j=1}^n \frac{1}{\sqrt{2\pi}} e^{-x_j^2/2} = \left(\frac{1}{2\pi}\right)^{n/2} e^{-\mathbf{x}^* \mathbf{x}/2}.$$

Because the standard normal distribution has mean 0, variance 1, and characteristic function $e^{-s^2/2}$, it follows that \mathbf{X} has mean vector $\mathbf{0}$, variance matrix \mathbf{I} , and characteristic function

$$E(e^{i\mathbf{s}^* \mathbf{X}}) = \prod_{j=1}^n e^{-s_j^2/2} = e^{-\mathbf{s}^* \mathbf{s}/2}.$$

We now define any affine transformation $\mathbf{Y} = \mathbf{A}\mathbf{X} + \boldsymbol{\mu}$ of \mathbf{X} to be multivariate normal [218]. This definition has several practical consequences. First, it is clear that $E(\mathbf{Y}) = \boldsymbol{\mu}$ and $\text{Var}(\mathbf{Y}) = \mathbf{A} \text{Var}(\mathbf{X}) \mathbf{A}^* = \mathbf{A} \mathbf{A}^* = \boldsymbol{\Omega}$. Second, any affine transformation $\mathbf{B}\mathbf{Y} + \boldsymbol{\nu} = \mathbf{B}\mathbf{A}\mathbf{X} + \mathbf{B}\boldsymbol{\mu} + \boldsymbol{\nu}$ of \mathbf{Y} is also multivariate normal. Third, any subvector of \mathbf{Y} is multivariate normal. Fourth, the characteristic function of \mathbf{Y} is

$$E(e^{i\mathbf{s}^* \mathbf{Y}}) = e^{i\mathbf{s}^* \boldsymbol{\mu}} E(e^{i\mathbf{s}^* \mathbf{A}\mathbf{X}}) = e^{i\mathbf{s}^* \boldsymbol{\mu} - \mathbf{s}^* \mathbf{A} \mathbf{A}^* \mathbf{s}/2} = e^{i\mathbf{s}^* \boldsymbol{\mu} - \mathbf{s}^* \boldsymbol{\Omega} \mathbf{s}/2}.$$

This enumeration omits two more subtle issues. One is whether \mathbf{Y} possesses a density. Observe that \mathbf{Y} lives in an affine subspace of dimension equal to or less than the rank of \mathbf{A} . Thus, if \mathbf{Y} has m components, then $n \geq m$ must hold in order for \mathbf{Y} to possess a density. A second issue is the existence and nature of the conditional density of a set of components of \mathbf{Y} given the remaining components. We can clarify both of these issues by making canonical choices of \mathbf{X} and \mathbf{A} based on the classical QR decomposition of a matrix, which follows directly from the Gram-Schmidt orthogonalization procedure [48].

Assuming that $n \geq m$, we can write

$$\mathbf{A}^* = \mathbf{Q} \begin{pmatrix} \mathbf{R} \\ \mathbf{0} \end{pmatrix},$$

where \mathbf{Q} is an $n \times n$ orthogonal matrix and $\mathbf{R} = \mathbf{L}^*$ is an $m \times m$ upper-triangular matrix with nonnegative diagonal entries. (If $n = m$, we omit the zero matrix in the QR decomposition.) It follows that

$$\mathbf{A}\mathbf{X} = (\mathbf{L} \ \mathbf{0}^*) \mathbf{Q}^* \mathbf{X} = (\mathbf{L} \ \mathbf{0}^*) \mathbf{Z}.$$

In view of the usual change-of-variables formula for probability densities and the facts that the orthogonal matrix \mathbf{Q}^* preserves inner products and has determinant ± 1 , the random vector \mathbf{Z} has n independent standard normal components and serves as a substitute for \mathbf{X} . Not only is this true, but we can dispense with the last $n - m$ components of \mathbf{Z} because they are multiplied by the matrix $\mathbf{0}^*$. Thus, we can safely assume $n = m$ and calculate the density of $\mathbf{Y} = \mathbf{L}\mathbf{Z} + \boldsymbol{\mu}$ when \mathbf{L} is invertible. In this situation, $\boldsymbol{\Omega} = \mathbf{L}\mathbf{L}^*$ is termed the Cholesky decomposition, and the usual change-of-variables formula shows that \mathbf{Y} has density

$$\begin{aligned} f(\mathbf{y}) &= \left(\frac{1}{2\pi}\right)^{n/2} |\det \mathbf{L}^{-1}| e^{-(\mathbf{y}-\boldsymbol{\mu})^*(\mathbf{L}^{-1})^* \mathbf{L}^{-1}(\mathbf{y}-\boldsymbol{\mu})/2} \\ &= \left(\frac{1}{2\pi}\right)^{n/2} |\det \boldsymbol{\Omega}|^{-1/2} e^{-(\mathbf{y}-\boldsymbol{\mu})^* \boldsymbol{\Omega}^{-1}(\mathbf{y}-\boldsymbol{\mu})/2}, \end{aligned}$$

where $\boldsymbol{\Omega} = \mathbf{L}\mathbf{L}^*$ is the variance matrix of \mathbf{Y} .

To address the issue of conditional densities, consider the compatibly partitioned vectors $\mathbf{Y}^* = (\mathbf{Y}_1^*, \mathbf{Y}_2^*)$, $\mathbf{X}^* = (\mathbf{X}_1^*, \mathbf{X}_2^*)$, and $\boldsymbol{\mu}^* = (\boldsymbol{\mu}_1^*, \boldsymbol{\mu}_2^*)$ and matrices

$$\mathbf{L} = \begin{pmatrix} \mathbf{L}_{11} & \mathbf{0} \\ \mathbf{L}_{21} & \mathbf{L}_{22} \end{pmatrix}, \quad \boldsymbol{\Omega} = \begin{pmatrix} \boldsymbol{\Omega}_{11} & \boldsymbol{\Omega}_{12} \\ \boldsymbol{\Omega}_{21} & \boldsymbol{\Omega}_{22} \end{pmatrix}.$$

Now suppose that \mathbf{X} is standard normal, that $\mathbf{Y} = \mathbf{L}\mathbf{X} + \boldsymbol{\mu}$, and that \mathbf{L}_{11} has full rank. For $\mathbf{Y}_1 = \mathbf{y}_1$ fixed, the equation $\mathbf{y}_1 = \mathbf{L}_{11}\mathbf{X}_1 + \boldsymbol{\mu}_1$ shows that \mathbf{X}_1 is fixed at the value $\mathbf{x}_1 = \mathbf{L}_{11}^{-1}(\mathbf{y}_1 - \boldsymbol{\mu}_1)$. Because no restrictions apply to \mathbf{X}_2 , we have

$$\mathbf{Y}_2 = \mathbf{L}_{22}\mathbf{X}_2 + \mathbf{L}_{21}\mathbf{L}_{11}^{-1}(\mathbf{y}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_2.$$

Thus, \mathbf{Y}_2 given \mathbf{Y}_1 is normal with mean $\mathbf{L}_{21}\mathbf{L}_{11}^{-1}(\mathbf{y}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_2$ and variance $\mathbf{L}_{22}\mathbf{L}_{22}^*$. To express these in terms of the blocks of $\boldsymbol{\Omega} = \mathbf{L}\mathbf{L}^*$, observe that

$$\begin{aligned} \boldsymbol{\Omega}_{11} &= \mathbf{L}_{11}\mathbf{L}_{11}^* \\ \boldsymbol{\Omega}_{21} &= \mathbf{L}_{21}\mathbf{L}_{11}^* \\ \boldsymbol{\Omega}_{22} &= \mathbf{L}_{21}\mathbf{L}_{21}^* + \mathbf{L}_{22}\mathbf{L}_{22}^*. \end{aligned}$$

The first two of these equations imply that $\mathbf{L}_{21}\mathbf{L}_{11}^{-1} = \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}$. The last equation then gives

$$\begin{aligned} \mathbf{L}_{22}\mathbf{L}_{22}^* &= \boldsymbol{\Omega}_{22} - \mathbf{L}_{21}\mathbf{L}_{21}^* \\ &= \boldsymbol{\Omega}_{22} - \boldsymbol{\Omega}_{21}(\mathbf{L}_{11}^*)^{-1}\mathbf{L}_{11}^{-1}\boldsymbol{\Omega}_{12} \\ &= \boldsymbol{\Omega}_{22} - \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}. \end{aligned}$$

These calculations do not require that \mathbf{Y}_2 possess a density. In summary, the conditional distribution of \mathbf{Y}_2 given \mathbf{Y}_1 is normal with mean and variance

$$\begin{aligned} \mathbf{E}(\mathbf{Y}_2 | \mathbf{Y}_1) &= \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}(\mathbf{Y}_1 - \boldsymbol{\mu}_1) + \boldsymbol{\mu}_2 \\ \text{Var}(\mathbf{Y}_2 | \mathbf{Y}_1) &= \boldsymbol{\Omega}_{22} - \boldsymbol{\Omega}_{21}\boldsymbol{\Omega}_{11}^{-1}\boldsymbol{\Omega}_{12}. \end{aligned} \tag{A.1}$$

A.3 Polyhedral Sets

A polyhedral set S is the nonempty intersection of a finite number of half-spaces. Symbolically, S can be represented as

$$S = \{ \mathbf{x} \in \mathbb{R}^m : \mathbf{v}_i^* \mathbf{x} \leq c_i \text{ for } 1 \leq i \leq p \}. \tag{A.2}$$

As previously noted, S is closed and convex. If all $c_i = 0$, then S is said to be a polyhedral cone. We now prove a sequence of propositions that lay out the basic facts about polyhedral cones and sets. Consult Examples 2.4.1 and 14.3.7 for background material on convex cones.

Proposition A.3.1 *The polar cone of a polyhedral cone is a finitely generated cone and vice versa. In matrix notation, the cones $\{ \mathbf{x} : \mathbf{V}^* \mathbf{x} \leq \mathbf{0} \}$ and $\{ \mathbf{V} \mathbf{a} : \mathbf{a} \geq \mathbf{0} \}$ constitute a polar pair.*

Proof: Assume the polyhedral set (A.2) is a cone. Let us show that its polar cone is

$$T = \left\{ \mathbf{y} \in \mathbb{R}^m : \mathbf{y} = \sum_{i=1}^p a_i \mathbf{v}_i, a_i \geq 0 \text{ for all } i \right\}.$$

According to Example 2.4.1, the finitely generated cone T is closed and convex. In view of Example 14.3.7, it therefore suffices to prove that the polar cone T° coincides with S . But this follows from the simple observation that $\mathbf{x}^* \mathbf{v}_i \leq 0$ for all \mathbf{v}_i if and only if $\mathbf{x}^* (\sum_{i=1}^p a_i \mathbf{v}_i) \leq 0$ for all conical combinations of the \mathbf{v}_i . ■

Proposition A.3.2 *A cone is polyhedral if and only if it is finitely generated.*

Proof: Consider a cone C generated by the vectors $\mathbf{u}_1, \dots, \mathbf{u}_p$ in \mathbb{R}^n . Let us prove by induction on p that C is a polyhedral cone. Suppose $p = 1$ and $\mathbf{u}_1 = \mathbf{0}$. If $\mathbf{e}_1, \dots, \mathbf{e}_n$ is the standard basis of \mathbb{R}^n , then

$$C = \{ \mathbf{0} \} = \{ \mathbf{x} : \mathbf{e}_i^* \mathbf{x} \leq 0, (-\mathbf{e}_i)^* \mathbf{x} \leq 0, \text{ for all } i \}.$$

If $p = 1$ and $\mathbf{u}_1 \neq \mathbf{0}$, then choose an orthogonal basis $\mathbf{w}_1, \dots, \mathbf{w}_n$ with $\mathbf{w}_1 = -\mathbf{u}_1$. In this basis, we have

$$C = \{ \mathbf{x} : \mathbf{w}_1^* \mathbf{x} \leq 0, \mathbf{w}_i^* \mathbf{x} \leq 0, (-\mathbf{w}_i)^* \mathbf{x} \leq 0, \text{ for all } i \geq 2 \}.$$

Now assume $p > 1$, and let K be the cone generated by $\mathbf{u}_1, \dots, \mathbf{u}_{p-1}$. By the induction hypothesis, K has polyhedral representation

$$K = \{ \mathbf{x} : \mathbf{v}_i^* \mathbf{x} \leq 0, 1 \leq i \leq q \}.$$

Furthermore, the cone C generated by $\mathbf{u}_1, \dots, \mathbf{u}_p$ satisfies

$$\begin{aligned} C &= \{ \mathbf{x} : \mathbf{x} - t\mathbf{u}_p \in K \text{ for some } t \geq 0 \} \\ &= \{ \mathbf{x} : \mathbf{v}_i^* (\mathbf{x} - t\mathbf{u}_p) \leq 0 \text{ for all } i \text{ and some } t \geq 0 \} \\ &= \{ \mathbf{x} : \mathbf{v}_i^* \mathbf{x} \leq t\mathbf{v}_i^* \mathbf{u}_p \text{ for all } i \text{ and some } t \geq 0 \}. \end{aligned}$$

To represent C as a polyhedral cone, we eliminate the variable t by the Fourier–Motzkin maneuver. Accordingly, define the index sets

$$I_- = \{i : \mathbf{v}_i^* \mathbf{u}_p < 0\}, \quad I_0 = \{i : \mathbf{v}_i^* \mathbf{u}_p = 0\}, \quad I_+ = \{i : \mathbf{v}_i^* \mathbf{u}_p > 0\}.$$

In this notation $\mathbf{x} \in C$ if and only if there exists $t \geq 0$ with

$$\frac{\mathbf{v}_i^* \mathbf{x}}{\mathbf{v}_i^* \mathbf{u}_p} \leq t \leq \frac{\mathbf{v}_j^* \mathbf{x}}{\mathbf{v}_j^* \mathbf{u}_p}, \quad \mathbf{v}_k^* \mathbf{x} \leq 0$$

for all $i \in I_+$, $j \in I_-$, and $k \in I_0$. Evidently, an appropriate $t \geq 0$ exists if and only if

$$\max_{i \in I_+} \frac{\mathbf{v}_i^* \mathbf{x}}{\mathbf{v}_i^* \mathbf{u}_p} \leq \min_{j \in I_-} \frac{\mathbf{v}_j^* \mathbf{x}}{\mathbf{v}_j^* \mathbf{u}_p} \quad \text{and} \quad \min_{j \in I_-} \frac{\mathbf{v}_j^* \mathbf{x}}{\mathbf{v}_j^* \mathbf{u}_p} \geq 0.$$

Because $\mathbf{v}_j^* \mathbf{u}_p < 0$ for $j \in I_-$, the second of the last two inequalities is equivalent to $\mathbf{v}_j^* \mathbf{x} \leq 0$ for all $j \in I_-$. It follows that C can be represented as the polyhedral cone

$$C = \left\{ \mathbf{x} : \mathbf{v}_k^* \mathbf{x} \leq 0, \frac{\mathbf{v}_i^* \mathbf{x}}{\mathbf{v}_i^* \mathbf{u}_p} \leq \frac{\mathbf{v}_j^* \mathbf{x}}{\mathbf{v}_j^* \mathbf{u}_p}, k \in I_- \cup I_0, i \in I_+, j \in I_- \right\}$$

involving no mention of the scalar t .

Conversely, if C is a polyhedral cone, then Proposition A.3.1 implies that C° is a finitely generated cone. By the argument just given, C° is also a polyhedral cone. A second application of Proposition A.3.1 shows that $C^{\circ\circ} = C$ is a finitely generated cone. ■

Proposition A.3.3 (Minkowski–Weyl) *A nonempty set S is polyhedral if and only if it can be represented as*

$$S = \left\{ \mathbf{x} : \mathbf{x} = \sum_{i=1}^q a_i \mathbf{u}_i + \sum_{j=1}^r b_j \mathbf{w}_j, \sum_{j=1}^r b_j = 1, \text{ all } a_i \geq 0, b_j \geq 0 \right\}. \quad (\text{A.3})$$

In other words, S is the algebraic sum of a finitely generated cone and the convex hull of a finite set of points.

Proof: Consider the polyhedral set S appearing in equation (A.2), and define the polyhedral cone

$$\begin{aligned} T &= \{ (\mathbf{x}, t) : t \geq 0 \text{ and } \mathbf{v}_i^* \mathbf{x} \leq c_i t \text{ for } 1 \leq i \leq p \} \\ &= \{ (\mathbf{x}, t) : t \geq 0 \text{ and } \mathbf{v}_i^* \mathbf{x} - c_i t \leq 0 \text{ for } 1 \leq i \leq p \}. \end{aligned}$$

Proposition A.3.2 implies that T is a finitely generated cone with generators $(\mathbf{u}_i, 0)$ and $(\mathbf{w}_j, 1)$ for $1 \leq i \leq q$ and $1 \leq j \leq r$. The representation (A.3) is now a consequence of the fact that $S = \{\mathbf{x} : (\mathbf{x}, 1) \in T\}$.

For the converse, suppose S is a set of the form (A.3). Define T to be the cone generated by the vectors $(\mathbf{u}_i, 0)$ and $(\mathbf{w}_j, 1)$ for $1 \leq i \leq q$ and $1 \leq j \leq r$. Proposition A.3.2 identifies T as a polyhedral cone satisfying

$$T = \{(\mathbf{x}, t) : t \geq 0 \text{ and } \mathbf{v}_i^* \mathbf{x} \leq c_i t \text{ for } 1 \leq i \leq p\}$$

for appropriate vectors $\mathbf{v}_1, \dots, \mathbf{v}_p$ and corresponding scalars c_1, \dots, c_p . But this means that $S = \{\mathbf{x} : (\mathbf{x}, 1) \in T\}$ is also a polyhedral set. ■

Proposition A.3.4 *The collection of polyhedral sets enjoys the following closure properties:*

- (a) *The nonempty intersection of two polyhedral sets is polyhedral.*
- (b) *The inverse image of a polyhedral set under a linear transformation is polyhedral.*
- (c) *The Cartesian product of two polyhedral sets is polyhedral.*
- (d) *The image of a polyhedral set under a linear transformation is polyhedral.*
- (e) *The vector sum of two polyhedral sets is polyhedral.*

Proof: Assertion (a) is obvious. Consider the polyhedral set (A.2). If \mathbf{M} is a linear transformation from \mathbb{R}^n to \mathbb{R}^m , then the set equality

$$\begin{aligned} \mathbf{M}^{-1}(S) &= \{\mathbf{x} \in \mathbb{R}^n : \mathbf{v}_i^*(\mathbf{M}\mathbf{x}) \leq c_i \text{ for } 1 \leq i \leq p\} \\ &= \{\mathbf{x} \in \mathbb{R}^n : (\mathbf{M}^* \mathbf{v}_i)^* \mathbf{x} \leq c_i \text{ for } 1 \leq i \leq p\} \end{aligned}$$

proves assertion (b). To verify assertion (c), consider a second polyhedral set

$$T = \{\mathbf{y} \in \mathbb{R}^n : \mathbf{w}_j^* \mathbf{y} \leq d_j \text{ for } 1 \leq j \leq q\}.$$

The Cartesian product amounts to

$$S \times T = \{(\mathbf{x}, \mathbf{y}) : \mathbf{v}_i^* \mathbf{x} \leq c_i, \mathbf{w}_j^* \mathbf{y} \leq d_j \text{ for all possible } i \text{ and } j\}.$$

To prove assertion (d), consider the Minkowski–Weyl representation

$$S = \left\{ \mathbf{x} \in \mathbb{R}^m : \mathbf{x} = \sum_{i=1}^p c_i \mathbf{v}_i + \sum_{j=1}^q d_j \mathbf{w}_j \right\},$$

where the first sum ranges over all convex combinations and the second sum ranges over all conical combinations. If \mathbf{N} is a linear transformation from \mathbf{R}^m to \mathbf{R}^n , then

$$\mathbf{N}(S) = \left\{ \mathbf{y} \in \mathbf{R}^n : \mathbf{y} = \sum_{i=1}^p c_i \mathbf{N} \mathbf{v}_i + \sum_{j=1}^q d_j \mathbf{N} \mathbf{w}_j \right\}$$

is a Minkowski–Weyl representation of the image. Finally, to prove assertion (e), note that $S + T$ is the image of the Cartesian product $S \times T$ under the linear map $(\mathbf{v}, \mathbf{w}) \mapsto \mathbf{v} + \mathbf{w}$. Application of properties (c) and (d) complete the proof. ■

Proposition A.3.5 *Let $f(\mathbf{x})$ be a convex function with domain a polyhedral set S . If $\sup_{\mathbf{x} \in S} f(\mathbf{x}) < \infty$, then $f(\mathbf{x})$ attains its maximum over S at one of the points \mathbf{w}_j defining the convex hull part of S .*

Proof: Assume that S has Minkowski–Weyl representation (A.3) and that $M = \sup_{\mathbf{x} \in S} f(\mathbf{x})$. If $\mathbf{x} \in S$ and \mathbf{u}_i is one of the vectors defining the conical part of S , then $\mathbf{x} + a_i \mathbf{u}_i \in S$ for all $a_i \geq 0$. Furthermore, for $t \geq 1$ we have

$$\begin{aligned} f(\mathbf{x} + a_i \mathbf{u}_i) &\leq \left(1 - \frac{1}{t}\right) f(\mathbf{x}) + \frac{1}{t} f(\mathbf{x} + t a_i \mathbf{u}_i) \\ &\leq \left(1 - \frac{1}{t}\right) f(\mathbf{x}) + \frac{M}{t}. \end{aligned}$$

Sending t to ∞ shows that $f(\mathbf{x} + a_i \mathbf{u}_i) \leq f(\mathbf{x})$. Hence, we may confine our attention to those points in the representation (A.3) with all $a_i = 0$. With this understanding, Jensen’s inequality

$$f\left(\sum_{j=1}^r b_j \mathbf{w}_j\right) \leq \sum_{j=1}^r b_j f(\mathbf{w}_j) \leq \max\{f(\mathbf{w}_1), \dots, f(\mathbf{w}_r)\}$$

demonstrates that $f(\mathbf{x})$ attains its maximum over S at one of the points \mathbf{w}_j defining the convex hull part of S . ■

A.4 Birkhoff’s Theorem and Fan’s Inequality

Birkhoff’s theorem deals with the set Γ^n of $n \times n$ doubly stochastic matrices. Every matrix $\mathbf{M} = (m_{ij})$ in Γ^n has nonnegative entries and row and column sums equal to 1. The affine constraints defining Γ^n compel it to be a compact polyhedral set with Minkowski–Weyl representation

$$\Gamma^n = \left\{ \mathbf{x} : \mathbf{x} = \sum_{j=1}^r b_j \mathbf{w}_j, \sum_{j=1}^r b_j = 1, b_j \geq 0 \right\} \quad (\text{A.4})$$

lacking a conical part. See Proposition A.3.3. Compact polyhedral sets are called convex polytopes and are characterized by their extreme points.

A point \mathbf{x} in a convex set S is said to be extreme if it cannot be written as a nontrivial convex combination of two points from S . It turns out that the Minkowski–Weyl vectors \mathbf{w}_i in a convex polytope can be taken to be extreme points. Indeed, suppose \mathbf{w}_i can be represented as

$$\mathbf{w}_i = \alpha \sum_{j=1}^r a_j \mathbf{w}_j + (1 - \alpha) \sum_{j=1}^r b_j \mathbf{w}_j$$

with $\alpha \in (0, 1)$. Either $a_i < 1$ or $b_i < 1$; otherwise, the two points on the right of the equation coincide with \mathbf{w}_i . Subtracting $[\alpha a_i + (1 - \alpha)b_i]\mathbf{w}_i$ from both sides of the equation and rescaling give \mathbf{w}_i as a convex combination $\sum_{j \neq i} c_j \mathbf{w}_j$. In any convex combination \mathbf{v} of the vectors $\{\mathbf{w}_j\}_{j=1}^r$, one can replace \mathbf{w}_i by this convex combination and represent \mathbf{v} by a convex combination of the remaining vectors $\mathbf{w}_j \neq \mathbf{w}_i$. If any non-extreme points remain after deletion of \mathbf{w}_i , then this substitution and reduction process can be repeated. Ultimately it halts with a set of vectors \mathbf{w}_j composed entirely of extreme points. In fact, these are the only extreme points of the convex polytope.

Birkhoff's theorem identifies the permutation matrices as the extreme points of Γ^n . A permutation matrix $\mathbf{P} = (p_{ij})$ has entries drawn from the set $\{0, 1\}$. Each of its rows and columns has exactly one entry equal to 1. The permutation matrices do not exhaust Γ^n . For instance, the matrix $\frac{1}{n}\mathbf{1}\mathbf{1}^*$ belongs to Γ^n . For another example, take any orthogonal matrix $\mathbf{U} = (u_{ij})$ and form the matrix \mathbf{M} with entries $m_{ij} = u_{ij}^2$. This matrix resides in Γ^n as well.

Proposition A.4.1 (Birkhoff) *Every doubly stochastic matrix can be represented as a convex combination of permutation matrices.*

Proof: It suffices to prove that the permutation matrices are the extreme points of Γ^n . Suppose the permutation matrix \mathbf{P} satisfies

$$\mathbf{P} = \alpha \mathbf{Q} + (1 - \alpha)\mathbf{R}$$

for two doubly stochastic matrices \mathbf{Q} and \mathbf{R} and $\alpha \in (0, 1)$. If an entry p_{ij} equals 1, then the two corresponding entries q_{ij} and r_{ij} of \mathbf{Q} and \mathbf{R} must also equal 1. Likewise, if an entry p_{ij} equals 0, then the corresponding entries q_{ij} and r_{ij} of \mathbf{Q} and \mathbf{R} must also equal 0. Thus, both \mathbf{Q} and \mathbf{R} coincide with \mathbf{P} . As a consequence, \mathbf{P} is an extreme point.

Conversely, suppose $\mathbf{M} = (m_{ij})$ is an extreme point that is not a permutation matrix. Take any index pair (i_0, j_0) with $0 < m_{i_0 j_0} < 1$. Because every row sum equals 1, there is an index $j_1 \neq j_0$ with $0 < m_{i_0 j_1} < 1$. Similarly, there is index i_1 with $0 < m_{i_1 j_1} < 1$. This process of jumping

along a row and then along a column creates a path from index pair to index pair. Eventually, the path intersects itself. Take a closed circuit

$$(i_k, j_k) \rightarrow (i_{k+1}, j_k) \rightarrow \cdots \rightarrow (i_l, j_l) = (i_k, j_k)$$

or

$$(i_{k+1}, j_k) \rightarrow (i_{k+1}, j_{k+1}) \rightarrow \cdots \rightarrow (i_l, j_l) = (i_{k+1}, j_k)$$

and construct a matrix \mathbf{N} whose entries are 0 except for entries along the path. For these special entries alternate the values 1 and -1 . It is clear that this construction forces \mathbf{N} to have row and column sums equal to 0. Because the entries of \mathbf{M} along the path occur in the open interval $(0, 1)$, there exists a positive constant ϵ such that $\mathbf{A} = \mathbf{M} + \epsilon\mathbf{N}$ and $\mathbf{B} = \mathbf{M} - \epsilon\mathbf{N}$ are both doubly stochastic. The representation

$$\mathbf{M} = \frac{1}{2}(\mathbf{A} + \mathbf{B})$$

now demonstrates that \mathbf{M} is not an extreme point. Hence, only permutation matrices can be extreme points. ■

As a prelude to stating Fan’s inequality, we first prove a classic rearrangement theorem of Hardy, Littlewood, and Pólya [118]. Consider two increasing sequences $a_1 \leq a_2 \leq \cdots \leq a_n$ and $b_1 \leq b_2 \leq \cdots \leq b_n$. If σ is any permutation of $\{1, \dots, n\}$, then the theorem says

$$\sum_{i=1}^n a_i b_{\sigma(i)} \leq \sum_{i=1}^n a_i b_i. \tag{A.5}$$

To quote the celebrated trio:

The theorem becomes obvious if we interpret the a_i as distances along a rod to hooks and the b_i as weights suspended from the hooks. To get the maximum statical moment with respect to the end of the rod, we hang the heaviest weights on the hooks farthest from the end.

To prove the result, suppose the a_i are in ascending order, but the b_i are not. Then there are indices $j < k$ with $a_j \leq a_k$ and $b_j > b_k$. Because

$$a_j b_k + a_k b_j - (a_j b_j + a_k b_k) = (a_k - a_j)(b_j - b_k) \geq 0,$$

we can increase the sum by exchanging b_j and b_k . A finite number of such exchanges (transpositions) puts the b_i into ascending order.

Proposition A.4.2 (von Neumann–Fan) *Let \mathbf{A} and \mathbf{B} be $n \times n$ symmetric matrices with ordered eigenvalues $\lambda_1 \geq \cdots \geq \lambda_n$ and $\mu_1 \geq \cdots \geq \mu_n$. Then*

$$\text{tr}(\mathbf{AB}) \leq \lambda_1 \mu_1 + \cdots + \lambda_n \mu_n. \tag{A.6}$$

Equality holds in inequality (A.6) if and only if

$$\mathbf{A} = \mathbf{W}\mathbf{D}_A\mathbf{W}^* \quad \text{and} \quad \mathbf{B} = \mathbf{W}\mathbf{D}_B\mathbf{W}^*$$

are simultaneously diagonalizable by an orthogonal matrix \mathbf{W} and diagonal matrices \mathbf{D}_A and \mathbf{D}_B whose entries are ordered from largest to smallest.

Proof: There exists a pair of orthogonal matrices \mathbf{U} and \mathbf{V} with

$$\mathbf{A} = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^* \quad \text{and} \quad \mathbf{B} = \sum_{i=1}^n \mu_i \mathbf{v}_i \mathbf{v}_i^*.$$

It follows that

$$\text{tr}(\mathbf{A}\mathbf{B}) = \sum_{i=1}^n \sum_{j=1}^n \lambda_i \mu_j \text{tr}(\mathbf{u}_i \mathbf{u}_i^* \mathbf{v}_j \mathbf{v}_j^*).$$

The matrix \mathbf{C} with entries $c_{ij} = \text{tr}(\mathbf{u}_i \mathbf{u}_i^* \mathbf{v}_j \mathbf{v}_j^*)$ is doubly stochastic because its entries $c_{ij} = (\mathbf{u}_i^* \mathbf{v}_j)^2$ are nonnegative and the column sums satisfy

$$\sum_{i=1}^n \text{tr}(\mathbf{u}_i \mathbf{u}_i^* \mathbf{v}_j \mathbf{v}_j^*) = \text{tr}\left(\sum_{i=1}^n \mathbf{u}_i \mathbf{u}_i^* \mathbf{v}_j \mathbf{v}_j^*\right) = \text{tr}(\mathbf{I}_n \mathbf{v}_j \mathbf{v}_j^*) = \mathbf{v}_j^* \mathbf{v}_j = 1.$$

Virtually the same argument shows that the row sums equal 1. Proposition A.3.3 therefore implies that \mathbf{C} is a convex combination $\sum_k \alpha_k \mathbf{P}_k$ of permutation matrices. This representation gives

$$\text{tr}(\mathbf{A}\mathbf{B}) = \boldsymbol{\lambda}^* \mathbf{C} \boldsymbol{\mu} = \sum_k \alpha_k \boldsymbol{\lambda}^* \mathbf{P}_k \boldsymbol{\mu} \leq \sum_k \alpha_k \sum_{i=1}^n \lambda_i \mu_i = \sum_{i=1}^n \lambda_i \mu_i$$

in view of the Hardy–Littlewood–Pólya rearrangement inequality (A.5).

Equality holds in inequality (A.6) under the stated conditions because

$$\text{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{W}\mathbf{D}_B\mathbf{W}^*) = \text{tr}(\mathbf{D}_A\mathbf{D}_B).$$

Here we apply the cyclic permutation property of the trace function and the identity $\mathbf{W}^*\mathbf{W} = \mathbf{I}_n$.

Conversely, suppose inequality (A.6) is an equality. Following Theobald [254], let $\mathbf{E} = \mathbf{A} + \mathbf{B}$ have ordered spectral decomposition $\mathbf{E} = \mathbf{W}\mathbf{D}_E\mathbf{W}^*$ with ρ_i the i th diagonal entry of \mathbf{D}_E . We now show that \mathbf{A} and \mathbf{B} have ordered spectral decompositions $\mathbf{W}\mathbf{D}_A\mathbf{W}^*$ and $\mathbf{W}\mathbf{D}_B\mathbf{W}^*$, respectively. The first half of the proposition and the hypothesis imply

$$\begin{aligned} \text{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{B}) &\leq \sum_{i=1}^n \lambda_i \mu_i = \text{tr}(\mathbf{A}\mathbf{B}) \\ \text{tr}(\mathbf{A}\mathbf{E}) &\leq \sum_{i=1}^n \lambda_i \rho_i = \text{tr}(\mathbf{D}_A\mathbf{D}_E). \end{aligned}$$

It follows that

$$\begin{aligned} \operatorname{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{A}) &= \operatorname{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{E}) - \operatorname{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{B}) \\ &\geq \operatorname{tr}(\mathbf{D}_A\mathbf{W}^*\mathbf{E}\mathbf{W}) - \operatorname{tr}(\mathbf{A}\mathbf{B}) \\ &= \operatorname{tr}(\mathbf{D}_A\mathbf{D}_E) - \operatorname{tr}(\mathbf{A}\mathbf{B}) \\ &\geq \operatorname{tr}(\mathbf{A}\mathbf{E}) - \operatorname{tr}(\mathbf{A}\mathbf{B}) \\ &= \operatorname{tr}(\mathbf{A}^2). \end{aligned}$$

Since $\operatorname{tr}(\mathbf{A}^2) = \|\mathbf{A}\|_F^2$ and $\|\mathbf{W}\mathbf{D}_A\mathbf{W}^*\|_F = \|\mathbf{D}_A\|_F = \|\mathbf{A}\|_F$, the Cauchy-Schwarz inequality for the Frobenius inner product now gives

$$\|\mathbf{A}\|_F^2 = \|\mathbf{W}\mathbf{D}_A\mathbf{W}^*\|_F\|\mathbf{A}\|_F \geq \operatorname{tr}(\mathbf{W}\mathbf{D}_A\mathbf{W}^*\mathbf{A}) \geq \|\mathbf{A}\|_F^2.$$

Because equality must hold in this inequality, the standard necessary condition for equality in the Cauchy-Schwarz inequality forces $\mathbf{W}\mathbf{D}_A\mathbf{W}^*$ to equal $c\mathbf{A}$ for some constant c . In fact, $c = 1$ because $\mathbf{W}\mathbf{D}_A\mathbf{W}^*$ has the same norm as \mathbf{A} . A similar argument implies that $\mathbf{B} = \mathbf{W}\mathbf{D}_B\mathbf{W}^*$. ■

Another proof of the sufficiency half of Proposition A.4.2 is possibly more illuminating [172]. Again let \mathbf{A} and \mathbf{B} be $n \times n$ symmetric matrices with eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$ and $\mu_1 \geq \dots \geq \mu_n$. One can easily show that the set of symmetric matrices with the same eigenvalues as \mathbf{A} is compact. Therefore the continuous function $\mathbf{M} \mapsto \operatorname{tr}(\mathbf{M}\mathbf{B})$ achieves its maximum over the set at some matrix \mathbf{A}_{max} . Now take any skew symmetric matrix \mathbf{C} and consider the one-parameter family of matrices $\mathbf{A}(t) = e^{t\mathbf{C}}\mathbf{A}_{max}e^{-t\mathbf{C}}$ similar to \mathbf{A}_{max} . Since $\mathbf{C}^* = -\mathbf{C}$, the matrix exponential $e^{t\mathbf{C}}$ is orthogonal with transpose $e^{-t\mathbf{C}}$. The optimality of \mathbf{A}_{max} implies that

$$\begin{aligned} \frac{d}{dt} \operatorname{tr}[\mathbf{A}(t)\mathbf{B}] \Big|_{t=0} &= \operatorname{tr}\{[\mathbf{C}\mathbf{A}_{max} - \mathbf{A}_{max}\mathbf{C}]\mathbf{B}\} \\ &= \operatorname{tr}[\mathbf{C}(\mathbf{A}_{max}\mathbf{B} - \mathbf{B}\mathbf{A}_{max})] \end{aligned}$$

vanishes. This suggests taking \mathbf{C} equal to the skew symmetric commutator matrix $\mathbf{A}_{max}\mathbf{B} - \mathbf{B}\mathbf{A}_{max}$. It then follows that

$$0 = \operatorname{tr}(\mathbf{C}\mathbf{C}) = -\operatorname{tr}(\mathbf{C}\mathbf{C}^*) = -\|\mathbf{C}\|_F^2.$$

In other words \mathbf{C} vanishes, and \mathbf{A}_{max} and \mathbf{B} commute. Commuting symmetric matrices can be simultaneously diagonalized by a common orthogonal matrix. Hence,

$$\operatorname{tr}(\mathbf{A}\mathbf{B}) \leq \operatorname{tr}(\mathbf{A}_{max}\mathbf{B}) = \sum_{i=1}^n \lambda_{\sigma(i)}\mu_i$$

for some permutation σ . Application of inequality (A.5) finishes the proof. The next proposition summarizes the foregoing discussion.

Proposition A.4.3 Consider an $n \times n$ symmetric matrix \mathbf{B} with ordered spectral decomposition $\mathbf{U} \operatorname{diag}(\boldsymbol{\mu})\mathbf{U}^*$. On the compact set of $n \times n$ symmetric matrices with ordered eigenvalues $\lambda_1 \geq \dots \geq \lambda_n$, the linear function

$$\mathbf{A} \mapsto \operatorname{tr}(\mathbf{A}\mathbf{B})$$

attains its maximum of $\sum_{i=1}^n \lambda_i \mu_i$ at the point $\mathbf{A} = \mathbf{U} \operatorname{diag}(\boldsymbol{\lambda})\mathbf{U}^*$.

A.5 Singular Value Decomposition

In many statistical applications involving large data sets, statisticians are confronted with a large $m \times n$ matrix $\mathbf{X} = (x_{ij})$ that encodes n features on each of m objects. For instance, in gene microarray studies x_{ij} represents the expression level of the i th gene under the j th experimental condition [190]. In information retrieval, x_{ij} represents the frequency of the j th word or term in the i th document [11]. The singular value decomposition (svd) captures the structure of such matrices. In many applications there are alternatives to the svd, but these are seldom as informative. From the huge literature on the svd, the books [64, 105, 107, 136, 137, 232, 249, 260] are especially recommended.

The spectral theorem for symmetric matrices discussed in Example 1.4.3 states that an $m \times m$ symmetric matrix \mathbf{M} can be written as $\mathbf{M} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^*$ for an orthogonal matrix \mathbf{U} and a diagonal matrix $\boldsymbol{\Lambda}$ with diagonal entries λ_i . If \mathbf{U} has columns $\mathbf{u}_1, \dots, \mathbf{u}_m$, then the matrix product $\mathbf{M} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^*$ unfolds into the sum of outer products

$$\mathbf{M} = \sum_{j=1}^m \lambda_j \mathbf{u}_j \mathbf{u}_j^*.$$

When $\lambda_j \neq 0$ for $j \leq k$ and $\lambda_j = 0$ for $j > k$, \mathbf{M} has rank k and only the first k terms of the sum are relevant. The svd seeks to generalize the spectral theorem to nonsymmetric matrices. In this case there are two orthonormal sets of vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$ and $\mathbf{v}_1, \dots, \mathbf{v}_k$ instead of one, and we write

$$\mathbf{M} = \sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^* = \mathbf{U}\boldsymbol{\Sigma}\mathbf{V}^* \tag{A.7}$$

for matrices \mathbf{U} and \mathbf{V} with orthonormal columns $\mathbf{u}_1, \dots, \mathbf{u}_k$ and $\mathbf{v}_1, \dots, \mathbf{v}_k$, respectively. If $\sigma_j < 0$, one can exchange $-\mathbf{u}_j$ for \mathbf{u}_j and $-\mathbf{v}_j$ for \mathbf{v}_j . Hence, it is possible to take the σ_j to be nonnegative in the representation (A.7).

For some purposes, it is better to fill out the matrices \mathbf{U} and \mathbf{V} to full orthogonal matrices. If \mathbf{M} is $m \times n$, then \mathbf{U} is viewed as $m \times m$, $\boldsymbol{\Sigma}$ as

$m \times n$, and \mathbf{V} as $n \times n$. The svd then becomes

$$\mathbf{M} = (\mathbf{u}_1 \dots \mathbf{u}_k \mathbf{u}_{k+1} \dots \mathbf{u}_m) \begin{pmatrix} \sigma_1^2 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & \sigma_k^2 & 0 & \dots & 0 \\ 0 & \dots & 0 & 0 & \dots & 0 \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ 0 & \dots & 0 & 0 & \dots & 0 \end{pmatrix} \begin{pmatrix} \mathbf{v}_1^* \\ \vdots \\ \mathbf{v}_k^* \\ \mathbf{v}_{k+1}^* \\ \vdots \\ \mathbf{v}_n^* \end{pmatrix},$$

assuming $k < \min\{m, n\}$. The scalars $\sigma_1, \dots, \sigma_k$ are said to be singular values and conventionally are listed in decreasing order. The vectors $\mathbf{u}_1, \dots, \mathbf{u}_k$ are known as left singular vectors and the vectors $\mathbf{v}_1, \dots, \mathbf{v}_k$ as right singular vectors.

To prove that the svd of an $m \times n$ matrix \mathbf{M} exists, consider the symmetric matrix product $\mathbf{M}^* \mathbf{M}$. Let $\mathbf{M}^* \mathbf{M}$ have spectral decomposition $\mathbf{V} \mathbf{\Omega} \mathbf{V}^*$ with the eigenvalues σ_i^2 arranged from greatest to least along the diagonal of $\mathbf{\Omega}$. The calculation

$$(\mathbf{M} \mathbf{v}_i)^* \mathbf{M} \mathbf{v}_j = \mathbf{v}_i^* \mathbf{M}^* \mathbf{M} \mathbf{v}_j = \sigma_j^2 \mathbf{v}_i^* \mathbf{v}_j = \begin{cases} 0 & i \neq j \\ \sigma_i^2 & i = j \end{cases}$$

shows that the vectors $\mathbf{M} \mathbf{v}_i$ are orthogonal. Furthermore, when $\sigma_i^2 > 0$, the normalized vector $\mathbf{u}_i = \sigma_i^{-1} \mathbf{M} \mathbf{v}_i$ is a unit vector. If we suppose that $\sigma_k > 0$ but $\sigma_{k+1} = 0$, then the representation (A.7) is valid because

$$\left(\sum_{j=1}^k \sigma_j \mathbf{u}_j \mathbf{v}_j^* \right) \mathbf{v}_i = \mathbf{M} \mathbf{v}_i$$

for all vectors in the orthonormal basis $\{\mathbf{v}_i\}_{i=1}^n$.

Fan's inequality generalizes to nonsymmetric matrices provided one substitutes ordered singular values for ordered eigenvalues. Let us first reduce the problem to square matrices. If \mathbf{M} is $m \times n$ with $n < m$, then the svd $\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$ translates into the svd

$$(\mathbf{M} \ \mathbf{0}) = \mathbf{U} (\mathbf{\Sigma} \ \mathbf{0}) \begin{pmatrix} \mathbf{V}^* & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{m-n} \end{pmatrix}$$

preserving the nontrivial singular values. Furthermore, the trace of two such expanded matrices satisfies

$$\text{tr} \left[(\mathbf{M} \ \mathbf{0})^* (\mathbf{N} \ \mathbf{0}) \right] = \text{tr}(\mathbf{M}^* \mathbf{N}).$$

A similar representation applies when $m < n$.

We now reduce the problem to symmetric matrices. Suppose \mathbf{M} is an $m \times m$ square matrix with svd $\mathbf{M} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$. The matrix

$$\mathbf{K} = \begin{pmatrix} \mathbf{0} & \mathbf{M} \\ \mathbf{M}^* & \mathbf{0} \end{pmatrix}$$

is a $2m \times 2m$ symmetric matrix, and the $2m$ vectors

$$\mathbf{w}_i = \frac{1}{\sqrt{2}} \begin{pmatrix} \mathbf{u}_i \\ \mathbf{v}_i \end{pmatrix} \quad \text{and} \quad \mathbf{z}_i = \frac{1}{\sqrt{2}} \begin{pmatrix} -\mathbf{u}_i \\ \mathbf{v}_i \end{pmatrix}$$

are orthogonal unit vectors. Because $M^*\mathbf{u}_i = \sigma_i\mathbf{v}_i$ and $M\mathbf{v}_i = \sigma_i\mathbf{u}_i$, the vector \mathbf{w}_i is an eigenvector of \mathbf{K} with eigenvalue σ_i , and the vector \mathbf{z}_i is an eigenvector of \mathbf{K} with eigenvalue $-\sigma_i$. Hence, we have found the spectral decomposition of \mathbf{K} .

Now let \mathbf{N} be a second $m \times m$ matrix with svd $P\Omega Q^*$ and a similar expansion to a $2m \times 2m$ symmetric matrix \mathbf{L} . Fan's inequality for symmetric matrices gives

$$\begin{aligned} 2 \operatorname{tr}(\mathbf{M}^*\mathbf{N}) &= \operatorname{tr}(\mathbf{M}^*\mathbf{N}) + \operatorname{tr}(\mathbf{N}\mathbf{M}^*) \\ &= \operatorname{tr}(\mathbf{M}^*\mathbf{N}) + \operatorname{tr}(\mathbf{M}\mathbf{N}^*) \\ &= \operatorname{tr}(\mathbf{K}^*\mathbf{L}) \\ &\leq \sum_{i=1}^m \sigma_i \omega_i + \sum_{i=1}^m (-\sigma_i)(-\omega_i) \\ &= 2 \sum_{i=1}^m \sigma_i \omega_i. \end{aligned}$$

Equality occurs in this inequality if and only if $\mathbf{P} = \mathbf{U}$ and $\mathbf{Q} = \mathbf{V}$. Thus, Fan's inequality extends to arbitrary matrices if we substitute ordered singular values for ordered eigenvalues.

A.6 Hadamard Semidifferentials

A function $f(\mathbf{y})$ mapping an open set U of \mathbb{R}^p into \mathbb{R}^q is said to be Hadamard semidifferentiable at $\mathbf{x} \in U$ if for every vector \mathbf{v} the uniform limit

$$\lim_{t \downarrow 0, \mathbf{w} \rightarrow \mathbf{v}} \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t} = d_{\mathbf{v}}f(\mathbf{x})$$

exists [63]. Taking \mathbf{w} identically equal to \mathbf{v} shows that $d_{\mathbf{v}}f(\mathbf{x})$ coincides with the forward directional derivative of $f(\mathbf{y})$ at \mathbf{x} in the direction \mathbf{v} . Some authors equate semidifferentiability at \mathbf{x} to the existence of all possible forward directional derivatives. Hadamard's definition is more restrictive and yields a richer theory. For the sake of brevity, we will omit the prefix Hadamard in discussing semidifferentials. It is also convenient to restate the definition in terms of sequences. Thus, semidifferentiability requires the limit

$$\lim_{n \rightarrow \infty} \frac{f(\mathbf{x} + t_n \mathbf{w}_n) - f(\mathbf{x})}{t_n} = d_{\mathbf{v}}f(\mathbf{x})$$

to exist and to be independent of the particular sequences $t_n \downarrow 0$ and $\mathbf{w}_n \rightarrow \mathbf{v}$. The relation between differentiability and semidifferentiability is spelled out in our first proposition.

Proposition A.6.1 *A function $f(\mathbf{y})$ differentiable at \mathbf{x} is also semidifferentiable at \mathbf{x} . If $f(\mathbf{y})$ is semidifferentiable at \mathbf{x} , and the map $\mathbf{v} \mapsto d_{\mathbf{v}}f(\mathbf{x})$ is linear, then $f(\mathbf{y})$ is differentiable at \mathbf{x} .*

Proof: If $f(\mathbf{y})$ differentiable at \mathbf{x} , then

$$f(\mathbf{y}) - f(\mathbf{x}) = df(\mathbf{x})(\mathbf{y} - \mathbf{x}) + o(\|\mathbf{y} - \mathbf{x}\|)$$

as \mathbf{y} approaches \mathbf{x} . Choosing $\mathbf{y} = \mathbf{x} + t\mathbf{w}$ for $t > 0$ and $\|\mathbf{w} - \mathbf{v}\|$ small shows that the Hadamard difference quotient approaches $d_{\mathbf{v}}f(\mathbf{x}) = df(\mathbf{x})\mathbf{v}$. To prove the partial converse, define

$$g(\mathbf{u}) = \frac{f(\mathbf{x} + \mathbf{u}) - f(\mathbf{x}) - d_{\mathbf{u}}f(\mathbf{x})}{\|\mathbf{u}\|},$$

and set $c = \limsup_{\|\mathbf{u}\| \rightarrow 0} \|g(\mathbf{u})\|$. It suffices to prove that $c = 0$. Choose a sequence $\mathbf{u}_n \neq \mathbf{0}$ such that $t_n = \|\mathbf{u}_n\|$ converges to 0 and $\|g(\mathbf{u}_n)\|$ converges to c . Because the unit sphere is compact, some subsequence of the sequence $\mathbf{w}_n = \|\mathbf{u}_n\|^{-1}\mathbf{u}_n$ converges to a unit vector \mathbf{v} . Without loss of generality, take the subsequence to be the original sequence. It then follows that

$$\begin{aligned} \lim_{n \rightarrow \infty} g(\mathbf{u}_n) &= \lim_{n \rightarrow \infty} \frac{f(\mathbf{x} + t_n \mathbf{w}_n) - f(\mathbf{x}) - d_{t_n \mathbf{w}_n}f(\mathbf{x})}{t_n} \\ &= \lim_{n \rightarrow \infty} \left[\frac{f(\mathbf{x} + t_n \mathbf{w}_n) - f(\mathbf{x})}{t_n} - d_{\mathbf{w}_n}f(\mathbf{x}) \right] \\ &= d_{\mathbf{v}}f(\mathbf{x}) - d_{\mathbf{v}}f(\mathbf{x}). \end{aligned}$$

The second equality here invokes the homogeneity of the map $\mathbf{v} \mapsto d_{\mathbf{v}}f(\mathbf{x})$. The third equality invokes the continuity of the map, which is a consequence of linearity. This calculation proves that $c = 0$. ■

Our second proposition demonstrates the utility of Hadamard's definition of semidifferentiability.

Proposition A.6.2 *A function $f(\mathbf{y})$ semidifferentiable at \mathbf{x} is continuous at \mathbf{x} .*

Proof: Suppose \mathbf{x}_n tends to \mathbf{x} but $f(\mathbf{x}_n)$ does not tend to $f(\mathbf{x})$. Then there exists an $\epsilon > 0$ such that $\|f(\mathbf{x}_n) - f(\mathbf{x})\| \geq \epsilon$ for infinitely many n . Without loss of generality, we may assume that the entire sequence possesses this property. Now write

$$\mathbf{x}_n = \mathbf{x} + \mathbf{x}_n - \mathbf{x} = \mathbf{x} + t_n \frac{\mathbf{x}_n - \mathbf{x}}{\|\mathbf{x}_n - \mathbf{x}\|}$$

by taking $t_n = \|\mathbf{x}_n - \mathbf{x}\|$. Again some subsequence of the sequence

$$\mathbf{w}_n = \frac{\mathbf{x}_n - \mathbf{x}_n}{\|\mathbf{x}_n - \mathbf{x}\|}$$

of unit vectors converges to a unit vector \mathbf{v} . Passing to the subsequence where this occurs if necessary, we have

$$f(\mathbf{x} + t_n \mathbf{w}_n) - f(\mathbf{x}) = t_n d_{\mathbf{v}} f(\mathbf{x}) + o(t_n), \tag{A.8}$$

contrary to the assumption that $\|f(\mathbf{x}_n) - f(\mathbf{x})\| \geq \epsilon$ for all large n . ■

The approximate equality (A.8) has an important consequence in minimization of real-valued functions. Suppose we can find a direction \mathbf{v} with $d_{\mathbf{v}} f(\mathbf{x}) < 0$. Then \mathbf{v} is a descent direction from \mathbf{x} in the sense that $f(\mathbf{x} + t\mathbf{v}) < f(\mathbf{x})$ for all sufficiently small $t > 0$. Thus, back-tracking is bound to produce a decrease in $f(\mathbf{y})$.

Here is a simple test for establishing semidifferentiability.

Proposition A.6.3 *Suppose $f(\mathbf{y})$ is locally Lipschitz around \mathbf{x} and possesses all possible forward directional derivatives there. Then $f(\mathbf{y})$ is semidifferentiable at \mathbf{x} .*

Proof: Consider the expansion

$$\frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t} = \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x} + t\mathbf{v})}{t} + \frac{f(\mathbf{x} + t\mathbf{v}) - f(\mathbf{x})}{t}.$$

As $t \downarrow 0$, the second fraction on the right of this equation tends to the forward directional derivative of $f(\mathbf{y})$ at \mathbf{x} in the direction \mathbf{v} . If c is a Lipschitz constant for $f(\mathbf{y})$ in a neighborhood of \mathbf{x} , then the first fraction on the right of the equation is locally bounded in norm by $c\|\mathbf{w} - \mathbf{v}\|$, which can be made arbitrarily small by taking \mathbf{w} close to \mathbf{v} . ■

Example A.6.1 *Convex Functions*

Any convex function $f(\mathbf{y})$ is locally Lipschitz and possesses all possible forward directional derivatives at an interior point \mathbf{x} of its essential domain. Proposition A.6.3 implies that $f(\mathbf{y})$ is semidifferentiable at such points. ■

Example A.6.2 *Norms*

A norm $\|\mathbf{x}\|_{\dagger}$ on \mathbb{R}^p is convex and therefore semidifferentiable in \mathbf{x} . In fact, a norm is globally Lipschitz because

$$\|\|\mathbf{y}\|_{\dagger} - \|\mathbf{x}\|_{\dagger}\| \leq \|\mathbf{y} - \mathbf{x}\|_{\dagger} \leq c\|\mathbf{y} - \mathbf{x}\|$$

for some constant $c > 0$. At the origin the semidifferential reduces to

$$\lim_{t \downarrow 0, \mathbf{w} \rightarrow \mathbf{v}} \frac{\|t\mathbf{w}\|_{\dagger} - 0}{t} = \|\mathbf{v}\|_{\dagger}.$$

At other points the semidifferential is more complicated to calculate. ■

Semidifferentiable functions obey most of the classical rules of differentiation.

Proposition A.6.4 *Let $f(\mathbf{y})$ and $g(\mathbf{y})$ be two functions semidifferentiable at the point \mathbf{x} . Then the homogeneity, sum, product, inverse, and chain rules*

$$\begin{aligned}d_{\mathbf{v}}[cf(\mathbf{x})] &= cd_{\mathbf{v}}f(\mathbf{x}) \\d_{\mathbf{v}}[f(\mathbf{x}) + g(\mathbf{x})] &= d_{\mathbf{v}}f(\mathbf{x}) + d_{\mathbf{v}}g(\mathbf{x}) \\d_{\mathbf{v}}[f(\mathbf{x})g(\mathbf{x})] &= d_{\mathbf{v}}f(\mathbf{x})g(\mathbf{x}) + f(\mathbf{x})d_{\mathbf{v}}g(\mathbf{x}) \\d_{\mathbf{v}}f(\mathbf{x})^{-1} &= -f(\mathbf{x})^{-1}d_{\mathbf{v}}f(\mathbf{x})f(\mathbf{x})^{-1} \\d_{\mathbf{v}}f \circ g(\mathbf{x}) &= d_{d_{\mathbf{v}}g(\mathbf{x})}f[g(\mathbf{x})]\end{aligned}$$

are valid under the usual compatibility assumptions for vector and matrix-valued functions.

Proof: These claims follow directly from the definition of semidifferentiability. Consider for instance the quotient rule. We simply write

$$\frac{1}{t} \left[\frac{1}{f(\mathbf{x} + t\mathbf{w})} - \frac{1}{f(\mathbf{x})} \right] = -f(\mathbf{x} + t\mathbf{w})^{-1} \cdot \frac{f(\mathbf{x} + t\mathbf{w}) - f(\mathbf{x})}{t} \cdot f(\mathbf{x})^{-1}$$

and take limits, invoking the continuity of $f(\mathbf{y})$ at \mathbf{x} in the process. For the chain rule, set

$$\mathbf{u} = \frac{g(\mathbf{x} + t\mathbf{w}) - g(\mathbf{x})}{t}$$

and rewrite the defining difference quotient as

$$\frac{f[g(\mathbf{x} + t\mathbf{w})] - f[g(\mathbf{x})]}{t} = \frac{f[g(\mathbf{x}) + t\mathbf{u}] - f[g(\mathbf{x})]}{t}.$$

Since \mathbf{u} tends to $d_{\mathbf{v}}g(\mathbf{x})$, the limit $d_{d_{\mathbf{v}}g(\mathbf{x})}f[g(\mathbf{x})]$ emerges. ■

Example A.6.3 *Semidifferential of $\|\mathbf{x}\|_r$ for $1 \leq r < \infty$*

The sum rule and a brief calculation give

$$d_{\mathbf{v}}\|\mathbf{x}\|_1 = \sum_{i=1}^p d_{\mathbf{v}}|x_i| = \sum_{i=1}^p \begin{cases} v_i & x_i > 0 \\ |v_i| & x_i = 0 \\ -v_i & x_i < 0. \end{cases}$$

Application of the chain and rules shows that the norm

$$\|\mathbf{x}\|_r = \left(\sum_{i=1}^p |x_i|^r \right)^{1/r}$$

for $1 < r < \infty$ has semidifferential

$$\begin{aligned} d_{\mathbf{v}}\|\mathbf{x}\|_r &= \frac{1}{r} \left(\sum_{i=1}^p |x_i|^r \right)^{\frac{1}{r}-1} \sum_{i=1}^p r|x_i|^{r-1} d_{\mathbf{v}}|x_i| \\ &= \|\mathbf{x}\|_r^{1-r} \sum_{i=1}^p |x_i|^{r-1} \operatorname{sgn}(x_i) v_i. \end{aligned}$$

Note that the semidifferential $d_{\mathbf{v}}\|\mathbf{x}\|_r$ does not necessarily converge to the semidifferential $d_{\mathbf{v}}\|\mathbf{x}\|_1$ as r tends to 1. ■

Example A.6.4 *Differences of Convex Functions*

If $f(\mathbf{y})$ and $g(\mathbf{y})$ are convex, then the difference $h(\mathbf{y}) = f(\mathbf{y}) - g(\mathbf{y})$ may not be convex. However, it is semidifferentiable throughout the intersection of the interiors of the essential domains of $f(\mathbf{y})$ and $g(\mathbf{y})$. ■

More surprising than the classical rules are the maxima and minima rules of the next proposition.

Proposition A.6.5 *Assume the real-valued functions $f_1(\mathbf{y}), \dots, f_m(\mathbf{y})$ are semi-differentiable at the point \mathbf{x} . If*

$$I(\mathbf{x}) = \{i : f_i(\mathbf{x}) = \max_i f_i(\mathbf{x})\} \text{ and } J(\mathbf{x}) = \{i : f_i(\mathbf{x}) = \min_i f_i(\mathbf{x})\},$$

then $\max_i f_i(\mathbf{y})$ and $\min_i f_i(\mathbf{y})$ are semidifferentiable at \mathbf{x} and

$$d_{\mathbf{v}} \max_i f_i(\mathbf{x}) = \max_{i \in I(\mathbf{x})} d_{\mathbf{v}} f_i(\mathbf{x}) \text{ and } d_{\mathbf{v}} \min_i f_i(\mathbf{x}) = \min_{i \in J(\mathbf{x})} d_{\mathbf{v}} f_i(\mathbf{x}).$$

Proof: The general rules follow from the case $m = 2$ and induction. Consider the minima rule. If $f_1(\mathbf{x}) < f_2(\mathbf{x})$, then this inequality persists in a neighborhood of \mathbf{x} . Hence, the rule follows by taking the limit of the difference quotient

$$\frac{\min\{f_1(\mathbf{x} + t\mathbf{w}), f_2(\mathbf{x} + t\mathbf{w})\} - \min\{f_1(\mathbf{x}), f_2(\mathbf{x})\}}{t} = \frac{f_1(\mathbf{x} + t\mathbf{w}) - f_1(\mathbf{x})}{t}.$$

The case $f_2(\mathbf{x}) < f_1(\mathbf{x})$ is handled similarly. For the case $f_2(\mathbf{x}) = f_1(\mathbf{x})$, we have

$$\frac{\min_{1 \leq i \leq 2} f_i(\mathbf{x} + t\mathbf{w}) - \min_{1 \leq i \leq 2} f_i(\mathbf{x})}{t} = \min_{1 \leq i \leq 2} \frac{f_i(\mathbf{x} + t\mathbf{w}) - f_i(\mathbf{x})}{t}.$$

Taking limits again validates the rule. ■

Example A.6.5 *Semidifferential of $\|\mathbf{x}\|_\infty$*

The maxima rule implies

$$d_{\mathbf{v}}\|\mathbf{x}\|_\infty = \max_{i \in I(\mathbf{x})} d_{\mathbf{v}}|x_i| = \max_{i \in I(\mathbf{x})} \begin{cases} v_i & x_i > 0 \\ |v_i| & x_i = 0 \\ -v_i & x_i < 0 \end{cases}$$

for $I(\mathbf{x}) = \{i : |x_i| = \|\mathbf{x}\|_\infty\}$. ■

Proposition A.6.5 was generalized by Danskin [54]. Here is a convex version of Danskin’s result.

Proposition A.6.6 *Consider a continuous function $f(\mathbf{x}, \mathbf{y})$ of two variables $\mathbf{x} \in \mathbb{R}^p$ and $\mathbf{y} \in C$ for some compact set $C \subset \mathbb{R}^q$. Suppose $f(\mathbf{x}, \mathbf{y})$ is convex in \mathbf{x} for each fixed \mathbf{y} . Then the function $g(\mathbf{x}) = \sup_{\mathbf{y} \in C} f(\mathbf{x}, \mathbf{y})$ is convex and has semidifferential*

$$d_{\mathbf{v}}g(\mathbf{x}) = \sup_{\mathbf{y} \in S(\mathbf{x})} d_{\mathbf{v}}f(\mathbf{x}, \mathbf{y}),$$

where $S(\mathbf{x})$ denotes the solution set of $\mathbf{y} \in C$ satisfying $f(\mathbf{x}, \mathbf{y}) = g(\mathbf{x})$. Finally, if $S(\mathbf{x})$ consists of a single point \mathbf{y} and $\nabla_{\mathbf{x}}f(\mathbf{x}, \mathbf{y})$ exists, then $\nabla g(\mathbf{x})$ exists and equals $\nabla_{\mathbf{x}}f(\mathbf{x}, \mathbf{y})$.

Proof: The function $g(\mathbf{x})$ is finite by virtue of the continuity of $f(\mathbf{x}, \mathbf{y})$ in \mathbf{y} and the compactness of C . It is convex because convexity is preserved under suprema. It therefore possesses a semidifferential, whose value equals the forward directional derivative $d_{\mathbf{v}}g(\mathbf{x})$. Now select any $\mathbf{y}_n \in S(\mathbf{x} + t_n\mathbf{v})$ and any $\mathbf{y} \in S(\mathbf{x})$. The inequality

$$\begin{aligned} \frac{g(\mathbf{x} + t_n\mathbf{v}) - g(\mathbf{x})}{t_n} &= \frac{f(\mathbf{x} + t_n\mathbf{v}, \mathbf{y}_n) - f(\mathbf{x}, \mathbf{y})}{t_n} \\ &= \frac{f(\mathbf{x} + t_n\mathbf{v}, \mathbf{y}_n) - f(\mathbf{x} + t_n\mathbf{v}, \mathbf{y})}{t_n} \\ &\quad + \frac{f(\mathbf{x} + t_n\mathbf{v}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y})}{t_n} \\ &\geq \frac{f(\mathbf{x} + t_n\mathbf{v}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y})}{t_n} \end{aligned}$$

implies in the limit that $d_{\mathbf{v}}g(\mathbf{x}) \geq d_{\mathbf{v}}f(\mathbf{x}, \mathbf{y})$ and consequently that

$$d_{\mathbf{v}}g(\mathbf{x}) \geq \sup_{\mathbf{y} \in S(\mathbf{x})} d_{\mathbf{v}}f(\mathbf{x}, \mathbf{y}).$$

In the process this logic reveals that $\sup_{\mathbf{y} \in S(\mathbf{x})} d_{\mathbf{v}}f(\mathbf{x}, \mathbf{y})$ is finite.

To prove the reverse inequality, observe that the inequalities

$$d_{\mathbf{v}}g(\mathbf{x}) \leq \frac{g(\mathbf{x} + t_n\mathbf{v}) - g(\mathbf{x})}{t_n}$$

$$\begin{aligned} &\leq \frac{f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) - f(\mathbf{x}, \mathbf{y}_n)}{t_n} \\ &\leq d_{\mathbf{v}} f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) \end{aligned}$$

for $\mathbf{y}_n \in S(\mathbf{x} + t_n \mathbf{v})$ simply reflect the monotonicity relations between difference quotients and directional derivatives for a convex function discussed in Sect. 6.4. To complete the proof, it suffices to argue that

$$\limsup_{n \rightarrow \infty} d_{\mathbf{v}} f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) \leq d_{\mathbf{v}} f(\mathbf{x}, \mathbf{y})$$

for some point $\mathbf{y} \in S(\mathbf{x})$. Fortunately, we can identify \mathbf{y} as the limit of any convergent subsequence of the original sequence \mathbf{y}_n . Without loss of generality, assume that this subsequence coincides with the original sequence. Taking limits in the inequality $f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) \geq f(\mathbf{x} + t_n \mathbf{v}, \mathbf{z})$ implies that $f(\mathbf{x}, \mathbf{y}) \geq f(\mathbf{x}, \mathbf{z})$ for all \mathbf{z} ; hence, $\mathbf{y} \in S(\mathbf{x})$. Now for any $\epsilon > 0$, all sufficiently small $t > 0$ satisfy

$$\frac{f(\mathbf{x} + t\mathbf{v}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y})}{t} \leq d_{\mathbf{v}} f(\mathbf{x}, \mathbf{y}) + \frac{\epsilon}{2}.$$

Hence, for such a t and sufficiently large n , joint continuity and monotonicity imply

$$\begin{aligned} d_{\mathbf{v}} f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) &\leq \frac{f(\mathbf{x} + t_n \mathbf{v} + t\mathbf{v}, \mathbf{y}_n) - f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n)}{t} \\ &\leq \frac{f(\mathbf{x} + t\mathbf{v}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y})}{t} + \frac{\epsilon}{2} \\ &\leq d_{\mathbf{v}} f(\mathbf{x}, \mathbf{y}) + \epsilon \end{aligned}$$

Since ϵ can be taken arbitrarily small in the inequality

$$d_{\mathbf{v}} f(\mathbf{x} + t_n \mathbf{v}, \mathbf{y}_n) \leq d_{\mathbf{v}} f(\mathbf{x}, \mathbf{y}) + \epsilon,$$

this completes the derivation of the semidifferential. To verify the last claim of the proposition, note that $d_{\mathbf{v}} g(\mathbf{x}) = d_{\mathbf{v}} f(\mathbf{x}, \mathbf{y}) = \nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})^* \mathbf{v}$. ■

Danskin’s original argument dispenses with convexity and relies on the existence and continuity of the gradient $\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})$. For our purposes the convex version is more convenient.

Example A.6.6 *Orthogonally Invariant Matrix Norms*

Let $\|\mathbf{A}\|_{\dagger}$ be a matrix norm on $m \times n$ matrices. As pointed out in Example 14.3.6, every matrix norm has a dual norm $\|\mathbf{B}\|_{\star}$ in terms of which

$$\|\mathbf{A}\|_{\dagger} = \sup_{\|\mathbf{B}\|_{\star}=1} \text{tr}(\mathbf{B}^* \mathbf{A}). \tag{A.9}$$

The semidifferential of the linear map $\mathbf{A} \mapsto \text{tr}(\mathbf{B}^* \mathbf{A})$ is

$$d_{\mathbf{C}} \text{tr}(\mathbf{B}^* \mathbf{A}) = \text{tr}(\mathbf{B}^* \mathbf{C}),$$

and the set $\{\mathbf{B} : \|\mathbf{B}\|_\star = 1\}$ is compact. Thus, the hypotheses of Proposition A.6.6 are met. Furthermore, representation (A.9) and duality show that $\|\cdot\|_\dagger$ is orthogonally invariant if and only if $\|\cdot\|_\star$ is orthogonally invariant. For an orthogonally invariant pair, the restriction $\|\mathbf{B}\|_\star = 1$ can be re-expressed as $\|\mathbf{D}_\mathbf{B}\|_\star = 1$ using the diagonal matrix $\mathbf{D}_\mathbf{B}$ whose diagonal entries are the singular values of \mathbf{B} .

According to Fan’s inequality, $\text{tr}(\mathbf{B}^*\mathbf{A}) \leq \sum_i \beta_i \alpha_i$, where the α_i are the ordered singular values of \mathbf{A} and the β_i are the ordered singular values of \mathbf{B} . Equality is attained in Fan’s inequality if and only if \mathbf{A} and \mathbf{B} have ordered singular value decompositions $\mathbf{A} = \mathbf{U}\mathbf{D}_\mathbf{A}\mathbf{V}^*$ and $\mathbf{B} = \mathbf{U}\mathbf{D}_\mathbf{B}\mathbf{V}^*$ with shared singular vectors. Let $S_\mathbf{A}$ be the set of diagonal matrices $\mathbf{D}_\mathbf{B}$ with ordered diagonal entries β_i , $\|\mathbf{D}_\mathbf{B}\|_\star = 1$, and $\sum_i \beta_i \alpha_i = \|\mathbf{A}\|_\dagger$. If the matrix \mathbf{U} has columns \mathbf{u}_i and the matrix \mathbf{V} has columns \mathbf{v}_i , then it follows that

$$d_C\|\mathbf{A}\|_\dagger = \sup_{\substack{\mathbf{D}_\mathbf{B} \in S_\mathbf{A} \\ \mathbf{A} = \mathbf{U}\mathbf{D}_\mathbf{B}\mathbf{V}^*}} \text{tr}(\mathbf{V}\mathbf{D}_\mathbf{B}\mathbf{U}^*\mathbf{C}) = \sup_{\substack{\mathbf{D}_\mathbf{B} \in S_\mathbf{A} \\ \mathbf{A} = \mathbf{U}\mathbf{D}_\mathbf{B}\mathbf{V}^*}} \sum_i \beta_i \mathbf{u}_i^* \mathbf{C} \mathbf{v}_i,$$

where the suprema extend over all singular value decompositions of \mathbf{A} . When the singular values are distinct, the singular vectors of \mathbf{A} are unique up to sign. The singular values are always unique. As an example, consider the spectral norm $\|\mathbf{B}\| = \beta_1$ and its dual the nuclear norm $\|\mathbf{B}\|_\star = \sum_i \beta_i$. The forward directional derivatives $d_C\|\mathbf{A}\|$ and $d_C\|\mathbf{A}\|_\star$ amount to

$$\begin{aligned} d_C\|\mathbf{A}\| &= \sup_{\mathbf{A} = \mathbf{U}\mathbf{D}_\mathbf{A}\mathbf{V}^*} \mathbf{u}_1^* \mathbf{C} \mathbf{v}_1 \\ d_C\|\mathbf{A}\|_\star &= \sup_{\mathbf{A} = \mathbf{U}\mathbf{D}_\mathbf{A}\mathbf{V}^*} \left(\sum_{\alpha_i > 0} \mathbf{u}_i^* \mathbf{C} \mathbf{v}_i + \sum_{\alpha_i = 0} \sup_{\beta_i \in [0,1]} \beta_i \mathbf{u}_i^* \mathbf{C} \mathbf{v}_i \right). \end{aligned}$$

In the first case we take $\beta_1 = 1$ and all remaining $\beta_i = 0$, and in the second case we take $\beta_i = 1$ when $\alpha_i > 0$ and $\beta_i \in [0, 1]$ otherwise. These directional derivatives are consistent with the rule (14.5) and the subdifferentials found in Example 14.5.7 and Problem 38 of Chap. 14. ■

Example A.6.7 *Induced Matrix Norms*

Let $\|\mathbf{x}\|_a$ and $\|\mathbf{y}\|_b$ be vector norms on \mathbb{R}^m and \mathbb{R}^n , respectively. These induce a norm on $m \times n$ matrices \mathbf{M} via

$$\|\mathbf{M}\|_{a,b} = \sup_{\|\mathbf{x}\|_a = 1} \|\mathbf{M}\mathbf{x}\|_b.$$

The chain rule implies $d_N\|\mathbf{M}\mathbf{x}\|_b = d_{N\mathbf{x}}\|\mathbf{y}\|_b$ with $\mathbf{y} = \mathbf{M}\mathbf{x}$. Hence, Proposition A.6.6 gives

$$d_N\|\mathbf{M}\|_{a,b} = \sup_{\substack{\|\mathbf{x}\|_a = 1 \\ \|\mathbf{M}\mathbf{x}\|_b = \|\mathbf{M}\|_{a,b}}} d_{N\mathbf{x}}\|\mathbf{M}\mathbf{x}\|_b.$$

In the special case where the two vectors norms are both ℓ_1 norms, one has $\|\mathbf{M}\|_{1,1} = \max_j \sum_i |m_{ij}|$. Suppose this maximum is attained for a unique index $j = k$. Then the best vector \mathbf{x} is the standard unit vector \mathbf{e}_k , and one can show that

$$d_{\mathcal{N}}\|\mathbf{M}\|_{1,1} = \sum_i \begin{cases} n_{ik} & m_{ik} > 0 \\ |n_{ik}| & m_{ik} = 0 \\ -n_{ik} & m_{ik} < 0. \end{cases}$$

The reader might like to consider the case of two ℓ_∞ vector norms. ■

Example A.6.8 *Differentiability of a Fenchel Conjugate*

Let $f(\mathbf{x})$ be a strictly convex function satisfying the growth condition $\liminf_{\|\mathbf{x}\| \rightarrow \infty} \|\mathbf{x}\|^{-1} f(\mathbf{x}) = \infty$. The Fenchel conjugate

$$f^*(\mathbf{y}) = \sup_{\mathbf{x}} [\mathbf{y}^* \mathbf{x} - f(\mathbf{x})]$$

is finite for all \mathbf{y} , and the supremum is attained at a unique point \mathbf{x} . In general for any compact set S of points \mathbf{y} , the corresponding set C of optimal points \mathbf{x} is compact. If on the contrary C is unbounded, then there exist paired sequences \mathbf{y}_n and \mathbf{x}_n with $\lim_{n \rightarrow \infty} \|\mathbf{x}_n\| = \infty$. Given the boundedness of S and the inequalities

$$-f(\mathbf{0}) \leq f^*(\mathbf{y}_n) = \mathbf{y}_n^* \mathbf{x}_n - f(\mathbf{x}_n) \leq \|\mathbf{y}_n\| \|\mathbf{x}_n\| - f(\mathbf{x}_n),$$

this contradicts the growth condition. The reader can check that C is closed. Propositions A.6.6 and A.6.1 now imply that $d_{\mathcal{N}} f^*(\mathbf{y}) = \mathbf{x}^* \mathbf{v}$ for the optimal \mathbf{x} and that $f^*(\mathbf{y})$ is differentiable with $\nabla f^*(\mathbf{y}) = \mathbf{x}$. ■

Example A.6.9 *Distance to a Closed Convex Set C*

The distance function $\text{dist}(\mathbf{x}, C)$ is convex and locally constant at an interior point of C . All directional derivatives consequently vanish there. At an exterior point the identity

$$\nabla \text{dist}(\mathbf{x}, C)^2 = 2[\mathbf{x} - P_C(\mathbf{x})] \tag{A.10}$$

involving the projection operator $P_C(\mathbf{x})$ shows that $\text{dist}(\mathbf{x}, C)$ is differentiable. Indeed, because $\text{dist}(\mathbf{x}, C)^2 > 0$, the chain rule yields

$$\nabla \text{dist}(\mathbf{x}, C) = \nabla \sqrt{\text{dist}(\mathbf{x}, C)^2} = \frac{\mathbf{x} - P_C(\mathbf{x})}{\text{dist}(\mathbf{x}, C)}.$$

To prove formula (A.10), set

$$\Delta = \text{dist}(\mathbf{x} + \mathbf{y}, C)^2 - \text{dist}(\mathbf{x}, C)^2.$$

In view of the inequality $\text{dist}(\mathbf{x}, C)^2 \leq \|\mathbf{x} - P_C(\mathbf{x} + \mathbf{y})\|^2$, one can construct the lower bound

$$\begin{aligned} \Delta &\geq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x} + \mathbf{y})\|^2 - \|\mathbf{x} - P_C(\mathbf{x} + \mathbf{y})\|^2 \\ &= \|\mathbf{y}\|^2 + 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x} + \mathbf{y})] \\ &= \|\mathbf{y}\|^2 + 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x})] + 2\mathbf{y}^*[P_C(\mathbf{x}) - P_C(\mathbf{x} + \mathbf{y})] \quad (\text{A.11}) \\ &\geq \|\mathbf{y}\|^2 + 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x})] - 2\|\mathbf{y}\| \cdot \|P_C(\mathbf{x}) - P_C(\mathbf{x} + \mathbf{y})\| \\ &\geq \|\mathbf{y}\|^2 + 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x})] - 2\|\mathbf{y}\|^2. \end{aligned}$$

The penultimate inequality here is just the Cauchy–Schwartz inequality. The final inequality is a consequence of the non-expansiveness of the projection operator. The analogous inequality $\text{dist}(\mathbf{x} + \mathbf{y}, C)^2 \leq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x})\|^2$ gives the upper bound

$$\Delta \leq \|\mathbf{x} + \mathbf{y} - P_C(\mathbf{x})\|^2 - \|\mathbf{x} - P_C(\mathbf{x})\|^2 = \|\mathbf{y}\|^2 + 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x})]. \quad (\text{A.12})$$

The two bounds (A.11) and (A.12) together imply that

$$\Delta = 2\mathbf{y}^*[\mathbf{x} - P_C(\mathbf{x})] + o(\|\mathbf{y}\|)$$

and consequently that $\nabla \text{dist}(\mathbf{x}, C)^2 = 2[\mathbf{x} - P_C(\mathbf{x})]$ according to Fréchet’s definition of the differential.

In contrast differentiability of $\text{dist}(\mathbf{x}, C)$ at boundary points of C is not assured. To calculate $d_{\mathbf{v}} \text{dist}(\mathbf{x}, C)$, we first observe that $\text{dist}(\mathbf{x}, C)$ equals the composition of the functions $f(\mathbf{y}) = \|\mathbf{y}\|$ and $g(\mathbf{x}) = \mathbf{x} - P_C(\mathbf{x})$. Given that $g(\mathbf{x}) = \mathbf{0}$, the chain rule therefore implies that

$$d_{\mathbf{v}} \text{dist}(\mathbf{x}, C) = \|\mathbf{v} - d_{\mathbf{v}}P_C(\mathbf{x})\|.$$

Thus, we need to evaluate $d_{\mathbf{v}}P_C(\mathbf{x})$. Without loss of generality, assume that $\mathbf{x} = \mathbf{0}$, and define $T(\mathbf{x})$ to be the closure of the convex cone $\cup_{t>0} \frac{1}{t}C$. Formally, $T(\mathbf{x})$ is the tangent space of C at the point $\mathbf{0}$. We now demonstrate that $d_{\mathbf{v}}P_C(\mathbf{x}) = P_{T(\mathbf{x})}(\mathbf{v})$ is the projection of \mathbf{v} onto $T(\mathbf{x})$. The easily checked identity

$$\frac{P_C(\mathbf{0} + t\mathbf{v}) - P_C(\mathbf{0})}{t} = P_{t^{-1}C}(\mathbf{v})$$

is our point of departure. We must demonstrate that $P_{t^{-1}C}(\mathbf{v})$ converges to $P_{T(\mathbf{x})}(\mathbf{v})$. Because $\mathbf{0} \in C$, the sets $t^{-1}C$ increase as t decreases. Hence, the trajectory $P_{t^{-1}C}(\mathbf{v})$ remains bounded. Let \mathbf{x} be any cluster point of $\mathbf{x}(t) = P_{t^{-1}C}(\mathbf{v})$. The obtuse angle criterion along the converging sequence $\mathbf{x}(t_n)$ requires

$$[\mathbf{v} - \mathbf{x}(t_n)]^*[\mathbf{z} - \mathbf{x}(t_n)] \leq 0$$

for every $\mathbf{z} \in t_n^{-1}C$. In the limit the inequality

$$(\mathbf{v} - \mathbf{x})^*(\mathbf{z} - \mathbf{x}) \leq 0$$

holds for every $z \in \cup_{t>0} \frac{1}{t}C$ and therefore for every z in the closure of this set. The only point in this closed set that qualifies is $P_{T(\mathbf{x})}(\mathbf{v})$. Thus, all cluster points of $\mathbf{x}(t)$ reduce to $P_{T(\mathbf{x})}(\mathbf{v})$. ■

A.7 Problems

1. Show that a polyhedral set S is compact if and only if it is the convex hull of a finite set of points.
2. Prove that a polyhedral set S possesses at most finitely many extreme points. (Hint: In the representation (A.2), a constraint $\mathbf{v}_i^* \mathbf{x} \leq c_i$ is said to be active at \mathbf{x} whenever equality holds there. If two points have the same active constraints, then neither of them is extreme.)
3. Demonstrate that every compact convex set possesses at least one extreme point. (Hint: The point farthest from the origin is extreme.)
4. Let $S = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ be a set of points in \mathbb{R}^n . If the pair $\{\mathbf{x}_i, \mathbf{x}_j\}$ attains the maximum Euclidean distance between two points in S , then show that \mathbf{x}_i and \mathbf{x}_j are both extreme points of $\text{conv}(S)$.
5. Let C_1 and C_2 be two closed convex sets. Demonstrate that

$$\text{dist}(C_1, C_2) = \inf_{\mathbf{x} \in C_1, \mathbf{y} \in C_2} \|\mathbf{x} - \mathbf{y}\|$$

is attained whenever (a) $C_1 \cap C_2 \neq \emptyset$, (b) either C_1 or C_2 is compact, or (c) both C_1 and C_2 are polyhedral.

6. For two sequences $a_1 \leq a_2 \leq \dots \leq a_n$ and $b_1 \leq b_2 \leq \dots \leq b_n$, demonstrate that

$$\sum_{i=1}^n a_i b_{\sigma(i)} \geq \sum_{i=1}^n a_i b_{n-i+1}$$

for every permutation σ .

7. Based on the previous exercise and the proof of Proposition A.4.2, devise a lower bound on $\text{tr}(\mathbf{A}\mathbf{B})$ for symmetric matrices \mathbf{A} and \mathbf{B} .
8. Under the Frobenius inner product $\langle \mathbf{M}, \mathbf{N} \rangle = \text{tr}(\mathbf{M}^* \mathbf{N})$ on square matrices, show that the subspaces S and A of symmetric and skew-symmetric matrices are orthogonal complements. Here $\mathbf{M} \in A$ if and only if $\mathbf{M} = -\mathbf{M}^*$. Find the projection operators P_S and P_A .

9. Suppose \mathbf{A} and \mathbf{B} are two $n \times n$ symmetric matrices with ordered eigenvalues $\{\lambda_i\}_{i=1}^n$ and $\{\mu_i\}_{i=1}^n$. Prove that

$$\sum_{i=1}^n (\lambda_i - \mu_i)^2 \leq \|\mathbf{A} - \mathbf{B}\|_F^2.$$

(Hint: $\sum_{i=1}^n \lambda_i^2 = \|\mathbf{A}\|_F^2$ and similarly for \mathbf{B} .)

10. The function

$$f(\mathbf{x}) = \begin{cases} \frac{x_1^6}{(x_2 - x_1^2)^2 + x_1^8} & \mathbf{x} \neq \mathbf{0} \\ 0 & \mathbf{x} = \mathbf{0} \end{cases}$$

illustrates the difference between Hadamard semidifferentiability at a point and the mere existence of all forward directional derivatives at the point. Prove the following assertions:

- The forward directional derivative $d_{\mathbf{v}}f(\mathbf{0}) = 0$ for all \mathbf{v} . (Hint: Treat the cases $v_2 = 0$ and $v_2 \neq 0$ separately.)
 - $f(\mathbf{x})$ is discontinuous at $\mathbf{0}$. (Hint: Take the limit of $f(t, t^2)$ as $t \rightarrow 0$.)
 - The Hadamard semidifferential $d_{(1,0)}f(\mathbf{0})$ does not exist. (Hint: Contrast the convergence of the relevant difference quotient for the scalar sequence $t_n = \frac{1}{n} \rightarrow 0$ and the two vector sequences $\mathbf{w}_n = (1, 0)^*$ and $\mathbf{w}_n = (1, \frac{1}{n})^* \rightarrow (1, 0)^*$.)
11. Demonstrate that a real-valued semidifferentiable function $f(\mathbf{x})$ satisfies

$$d_{\mathbf{v}}|f(\mathbf{x})| = \begin{cases} d_{\mathbf{v}}f(\mathbf{x}) & f(\mathbf{x}) > 0 \\ |d_{\mathbf{v}}f(\mathbf{x})| & f(\mathbf{x}) = 0 \\ -d_{\mathbf{v}}f(\mathbf{x}) & f(\mathbf{x}) < 0 \end{cases}$$

and

$$d_{\mathbf{v}} \max\{f(\mathbf{x}), 0\} = \begin{cases} d_{\mathbf{v}}f(\mathbf{x}) & f(\mathbf{x}) > 0 \\ \max\{d_{\mathbf{v}}f(\mathbf{x}), 0\} & f(\mathbf{x}) = 0 \\ 0 & f(\mathbf{x}) < 0. \end{cases}$$

These formulas are pertinent in calculating the directional derivatives of the exact penalty function $\mathcal{E}_\rho(\mathbf{y})$ discussed in Sect. 16.3.

12. Suppose the function $f(\mathbf{x})$ is Lipschitz in a neighborhood of the point \mathbf{y} with Lipschitz constant c . Prove the inequality

$$\|d_{\mathbf{v}}f(\mathbf{y}) - d_{\mathbf{w}}f(\mathbf{y})\| \leq c\|\mathbf{v} - \mathbf{w}\|$$

assuming the indicated forward directional derivatives exist. The special case $\mathbf{w} = \mathbf{0}$ gives $\|d_{\mathbf{v}}f(\mathbf{x})\| \leq c\|\mathbf{v}\|$.

References

- [1] Acosta E, Delgado C (1994) Fréchet versus Carathéodory. *Am Math Mon* 101:332–338
- [2] Acton FS (1990) *Numerical methods that work*. Mathematical Association of America, Washington, DC
- [3] Anderson TW (2003) *An introduction to multivariate statistical analysis*, 3rd edn. Wiley, Hoboken
- [4] Armstrong RD, Kung MT (1978) Algorithm AS 132: least absolute value estimates for a simple linear regression problem. *Appl Stat* 27:363–366
- [5] Arthur D, Vassilvitskii S (2007) k-means++: the advantages of careful seeding. In: *2007 symposium on discrete algorithms (SODA)*. Society for Industrial and Applied Mathematics, Philadelphia, 2007
- [6] Barlow RE, Bartholomew DJ, Bremner JM, Brunk HD (1972) *Statistical inference under order restrictions; the theory and application of isotonic regression*. Wiley, New York
- [7] Bartle RG (1996) Return to the Riemann integral. *Am Math Mon* 103:625–632
- [8] Baum LE (1972) An inequality and associated maximization technique in statistical estimation for probabilistic functions of Markov processes. *Inequalities* 3:1–8

- [9] Bauschke HH, Lewis AS (2000) Dykstra's algorithm with Bregman projections: a convergence proof. *Optimization* 48:409–427
- [10] Beltrami EJ (1970) An algorithmic approach to nonlinear analysis and optimization. Academic, New York
- [11] Berry MW, Drmac Z, Jessup ER (1999) Matrices, vector spaces, and information retrieval. *SIAM Rev* 41:335–362
- [12] Bertsekas DP (1999) Nonlinear programming, 2nd edn. Athena Scientific, Belmont
- [13] Bertsekas DP (2009) Convex optimization theory. Athena Scientific, Belmont
- [14] Bishop YMM, Feinberg SE, Holland PW (1975) Discrete multivariate analysis: theory and practice. MIT, Cambridge
- [15] Bliss GA (1925) Calculus of variations. Mathematical Society of America, Washington, DC
- [16] Böhning D, Lindsay BG (1988) Monotonicity of quadratic approximation algorithms. *Ann Inst Stat Math* 40:641–663
- [17] Borwein JM, Lewis AS (2000) Convex analysis and nonlinear optimization: theory and examples. Springer, New York
- [18] Botsko MW, Gosser RA (1985) On the differentiability of functions of several variables. *Am Math Mon* 92:663–665
- [19] Boyd S, Vandenberghe L (2004) Convex optimization. Cambridge University Press, Cambridge
- [20] Boyd S, Kim SJ, Vandenberghe L, Hassibi A (2007) A tutorial on geometric programming. *Optim Eng* 8:67–127
- [21] Boyle JP, Dykstra RL (1985) A method for finding projections onto the intersection of convex sets in Hilbert space. In: *Advances in order restricted statistical inference. Lecture notes in statistics*. Springer, New York, pp 28–47
- [22] Bradley EL (1973) The equivalence of maximum likelihood and weighted least squares estimates in the exponential family. *J Am Stat Assoc* 68:199–200
- [23] Bradley RA, Terry ME (1952) Rank analysis of incomplete block designs. *Biometrika* 39:324–345
- [24] Bregman LM (1965) The method of successive projection for finding a common point of convex sets. *Sov Math Dokl* 6:688–692

- [25] Bregman LM (1967) The relaxation method of finding the common points of convex sets and its application to the solution of problems in convex programming. *USSR Comput Math Math Phys* 7:200–217
- [26] Bregman LM, Censor Y, Reich S (2000) Dykstra’s algorithm as the nonlinear extension of Bregman’s optimization method. *J Convex Anal* 6:319–333
- [27] Brent RP (1973) Some efficient algorithms for solving systems of nonlinear equations. *SIAM J Numer Anal* 10:327–344
- [28] Brezhneva OA, Tret’yakov AA, Wright SE (2010) A simple and elementary proof of the Karush-Kuhn-Tucker theorem for inequality-constrained optimization. *Optim Lett* 3:7–10
- [29] Bridger M, Stolzenberg G (1999) Uniform calculus and the law of bounded change. *Am Math Mon* 106:628–635
- [30] Brinkhuis J, Tikhomirov V (2005) *Optimization: insights and applications*. Princeton University Press, Princeton
- [31] Brophy JF, Smith PW (1988) Prototyping Karmarkar’s algorithm using MATH-PROTRAN. *IMSL Dir* 5:2–3
- [32] Broyden CG (1965) A class of methods for solving nonlinear simultaneous equations. *Math Comput* 19:577–593
- [33] Byrd RH, Nocedal J (1989) A tool for the analysis of quasi-Newton methods with application to unconstrained minimization. *SIAM J Numer Anal* 26:727–739
- [34] Byrne CL (2009) *A first course in optimization*. Department of Mathematical Sciences, University of Massachusetts Lowell, Lowell
- [35] Cai J-F, Candés EJ, Shen Z (2008) A singular value thresholding algorithm for matrix completion. *SIAM J Optim* 20:1956–1982
- [36] Candés EJ, Tao T (2007) The Danzig selector: statistical estimation when p is much larger than n . *Ann Stat* 35:2313–2351
- [37] Candés EJ, Tao T (2009) The power of convex relaxation: near-optimal matrix completion. *IEEE Trans Inform Theor* 56:2053–2080
- [38] Candés EJ, Romberg J, Tao T (2006) Stable signal recovery from incomplete and inaccurate measurements. *Comm Pure Appl Math* 59:1207–1223
- [39] Candés EJ, Wakin M, Boyd S (2007) Enhancing sparsity by reweighted ℓ_1 minimization. *J Fourier Anal Appl* 14:877–905

- [40] Carathéodory C (1954) Theory of functions of a complex variable, vol 1. Chelsea, New York
- [41] Censor Y, Zenios SA (1992) Proximal minimization with D-functions. *J Optim Theor Appl* 73:451–464
- [42] Censor Y, Chen W, Combettes PL, Davidi R, Herman GT (2012) On the effectiveness of projection methods for convex feasibility problems with linear inequality constraints. *Comput Optim Appl* 51:1065–1088
- [43] Charnes A, Frome EL, Yu PL (1976) The equivalence of generalized least squares and maximum likelihood in the exponential family. *J Am Stat Assoc* 71:169–171
- [44] Chen J, Tan X (2009) Inference for multivariate normal mixtures. *J Multivariate Anal* 100:1367–1383
- [45] Chen SS, Donoho DL, Saunders MA (1998) Atomic decomposition by basis pursuit. *SIAM J Sci Comput* 20:33–61
- [46] Cheney W (2001) Analysis for applied mathematics. Springer, New York
- [47] Choi SC, Wette R (1969) Maximum likelihood estimation of the parameters of the gamma distribution and their bias. *Technometrics* 11:683–690
- [48] Ciarlet PG (1989) Introduction to numerical linear algebra and optimization. Cambridge University Press, Cambridge
- [49] Claerbout J, Muir F (1973) Robust modeling with erratic data. *Geophysics* 38:826–844
- [50] Clarke CA, Price Evans DA, McConnell RB, Sheppard PM (1959) Secretion of blood group antigens and peptic ulcers. *Br Med J* 1:603–607
- [51] Conn AR, Gould NIM, Toint PL (1991) Convergence of quasi-Newton matrices generated by the symmetric rank one update. *Math Program* 50:177–195
- [52] Conte SD, deBoor C (1972) Elementary numerical analysis. McGraw-Hill, New York
- [53] Cox DR (1970) Analysis of binary data. Methuen, London
- [54] Danskin JM (1966) The theory of max-min, with applications. *SIAM J Appl Math* 14:641–664

- [55] Daubechies I, Defrise M, De Mol C (2004) An iterative thresholding algorithm for linear inverse problems with a sparsity constraint. *Comm Pure Appl Math* 57:1413–1457
- [56] Davidon WC (1959) Variable metric methods for minimization. AEC Research and Development Report ANL-5990, Argonne National Laboratory, Argonne
- [57] Davis JA, Smith TW (2008) General social surveys, 1972–2008 [machine-readable data file]. Roper Center for Public Opinion Research, University of Connecticut, Storrs
- [58] Debreu G (1952) Definite and semidefinite quadratic forms. *Econometrica* 20:295–300
- [59] de Leeuw J (1994) Block relaxation algorithms in statistics. In: Bock HH, Lenski W, Richter MM (eds) *Information systems and data analysis*. Springer, New York, pp 308–325
- [60] de Leeuw J (2006) Some majorization techniques. Preprint series, UCLA Department of Statistics.
- [61] de Leeuw J, Heiser WJ (1980) Multidimensional scaling with restrictions on the configuration. In: Krishnaiah PR (ed) *Multivariate analysis*, vol V. North-Holland, Amsterdam, pp 501–522
- [62] de Leeuw J, Lange K (2009) Sharp quadratic majorization in one dimension. *Comput Stat Data Anal* 53:2471–2484
- [63] Delfour MC (2012) *Introduction to optimization and semidifferential calculus*. SIAM, Philadelphia
- [64] Demmel J (1997) *Applied numerical linear algebra*. SIAM, Philadelphia
- [65] Dempster AP, Laird NM, Rubin DB (1977) Maximum likelihood from incomplete data via the EM algorithm (with discussion). *J Roy Stat Soc B* 39:1–38
- [66] Dennis JE Jr, Schnabel RB (1996) *Numerical methods for unconstrained optimization and nonlinear equations*. SIAM, Philadelphia
- [67] De Pierro AR (1993) On the relation between the ISRA and EM algorithm for positron emission tomography. *IEEE Trans Med Imag* 12:328–333
- [68] DePree JD, Swartz CW (1988) *Introduction to real analysis*. Wiley, Hoboken

- [69] de Souza PN, Silva J-N (2001) Berkeley problems in mathematics, 2nd edn. Springer, New York
- [70] Deutsch F (2001) Best approximation in inner product spaces. Springer, New York
- [71] Devijver PA (1985) Baum's forward-backward algorithm revisited. *Pattern Recogn Lett* 3:369–373
- [72] Ding C, Li T, Jordan MI (2010) Convex and semi-nonnegative matrix factorizations. *IEEE Trans Pattern Anal Mach Intell* 32:45–55
- [73] Dobson AJ (1990) An introduction to generalized linear models. Chapman & Hall, London
- [74] Donoho DL (2006) Compressed sensing. *IEEE Trans Inform Theor* 52:1289–1306
- [75] Donoho D, Johnstone I (1994) Ideal spatial adaptation by wavelet shrinkage. *Biometrika* 81:425–455
- [76] Duan J-C, Simonato J-G (1993) Multiplicity of solutions in maximum likelihood factor analysis. *J Stat Comput Simul* 47:37–47
- [77] Duchi J, Shalev-Shwartz S, Singer Y, Chandra T (2008) Efficient projections onto the l_1 -ball for learning in high dimensions. In: Proceedings of the 25th international conference on machine learning (ICML 2008). ACM, New York, pp 272–279
- [78] Durbin R, Eddy S, Krogh A, Mitchison G (1998) Biological sequence analysis: probabilistic models of proteins and nucleic acids. Cambridge University Press, Cambridge
- [79] Dykstra RL (1983) An algorithm for restricted least squares estimation. *J Am Stat Assoc* 78:837–842
- [80] Edgeworth FY (1887) On observations relating to several quantities. *Hermathena* 6:279–285
- [81] Edgeworth FY (1888) On a new method of reducing observations relating to several quantities. *Phil Mag* 25:184–191
- [82] Edwards CH Jr (1973) Advanced calculus of several variables. Academic, New York
- [83] Efron B, Hastie T, Johnstone I, Tibshirani R (2004) Least angle regression. *Ann Stat* 32:407–499
- [84] Ekeland I (1974) On the variational principle. *J Math Anal Appl* 47:324–353

- [85] Elsner L, Koltracht L, Neumann M (1992) Convergence of sequential and asynchronous nonlinear paracontractions. *Numer Math* 62:305–319
- [86] Everitt BS (1977) *The analysis of contingency tables*. Chapman & Hall, London
- [87] Fang S-C, Puthenpura S (1993) *Linear optimization and extensions: theory and algorithms*. Prentice-Hall, Englewood Cliffs
- [88] Fazel M, Hindi M, Boyd S (2003) Log-det heuristic for matrix rank minimization with applications to Hankel and Euclidean distance matrices. *Proc Am Contr Conf* 3:2156–2162
- [89] Feller W (1971) *An introduction to probability theory and its applications*, vol 2, 2nd edn. Wiley, Hoboken
- [90] Fessler JA, Clinthorne NH, Rogers WL (1993) On complete-data spaces for PET reconstruction algorithms. *IEEE Trans Nucl Sci* 40:1055–1061
- [91] Fiacco AV, McCormick GP (1968) *Nonlinear programming: sequential unconstrained minimization techniques*. Wiley, Hoboken
- [92] Fletcher R (2000) *Practical methods of optimization*, 2nd edn. Wiley, Hoboken
- [93] Fletcher R, Powell MJD (1963) A rapidly convergent descent method for minimization. *Comput J* 6:163–168
- [94] Fletcher R, Reeves CM (1964) Function minimization by conjugate gradients. *Comput J* 7:149–154
- [95] Flury B, Zoppè A (2000) Exercises in EM. *Am Stat* 54:207–209
- [96] Forsgren A, Gill PE, Wright MH (2002) Interior point methods for nonlinear optimization. *SIAM Rev* 44:523–597
- [97] Franklin J (1983) Mathematical methods of economics. *Am Math Mon* 90:229–244
- [98] Friedman J, Hastie T, Tibshirani R (2007) Pathwise coordinate optimization. *Ann Appl Stat* 1:302–332
- [99] Friedman J, Hastie T, Tibshirani R (2009) Regularized paths for generalized linear models via coordinate descent. Technical Report, Department of Statistics, Stanford University
- [100] Fu WJ (1998) Penalized regressions: the bridge versus the lasso. *J Comput Graph Stat* 7:397–416

- [101] Gabriel KR, Zamir S (1979) Lower rank approximation of matrices by least squares with any choice of weights. *Technometrics* 21:489–498
- [102] Gelfand IM, Fomin SV (1963) *Calculus of variations*. Prentice-Hall, Englewood Cliffs
- [103] Geman S, McClure D (1985) Bayesian image analysis: an application to single photon emission tomography. In: *Proceedings of the statistical computing section*. American Statistical Association, Washington, DC, pp 12–18
- [104] Gifi A (1990) *Nonlinear multivariate analysis*. Wiley, Hoboken
- [105] Gill PE, Murray W, Wright MH (1991) *Numerical linear algebra and optimization*, vol 1. Addison-Wesley, Redwood City
- [106] Goldstein T, Osher S (2009) The split Bregman method for ℓ_1 -regularized problems. *SIAM J Imag Sci* 2:323–343
- [107] Golub GH, Van Loan CF (1996) *Matrix computations*, 3rd edn. Johns Hopkins University Press, Baltimore
- [108] Gordon RA (1998) The use of tagged partitions in elementary real analysis. *Am Math Mon* 105:107–117
- [109] Gould NIM (2008) How good are projection methods for convex feasibility problems? *Comput Optim Appl* 40:1–12
- [110] Green PJ (1984) Iteratively reweighted least squares for maximum likelihood estimation and some robust and resistant alternatives (with discussion). *J Roy Stat Soc B* 46:149–192
- [111] Green PJ (1990) Bayesian reconstruction for emission tomography data using a modified EM algorithm. *IEEE Trans Med Imag* 9:84–94
- [112] Green PJ (1990) On use of the EM algorithm for penalized likelihood estimation. *J Roy Stat Soc B* 52:443–452
- [113] Grimmett GR, Stirzaker DR (1992) *Probability and random processes*, 2nd edn. Oxford University Press, Oxford
- [114] Groenen PJF, Nalbantov G, Bioch JC (2007) Nonlinear support vector machines through iterative majorization and I-splines. In: Lenz HJ, Decker R (eds) *Studies in classification, data analysis, and knowledge organization*. Springer, Heidelberg, pp 149–161
- [115] Guillemin V, Pollack A (1974) *Differential topology*. Prentice-Hall, Englewood Cliffs
- [116] Güler O (2010) *Foundations of optimization*. Springer, New York

- [117] Hämmerlin G, Hoffmann K-H (1991) Numerical mathematics. Springer, New York
- [118] Hardy GH, Littlewood JE, Pólya G (1952) Inequalities, 2nd edn. Cambridge University Press, Cambridge
- [119] Hastie T, Tibshirani R, Friedman J (2009) The elements of statistical learning: data mining, inference, and prediction, 2nd edn. Springer, New York
- [120] He L, Marquina A, Osher S (2005) Blind deconvolution using TV regularization and Bregman iteration. *Int J Imag Syst Technol* 15, 74–83
- [121] Heiser WJ (1987) Correspondence analysis with least absolute residuals. *Comput Stat Data Anal* 5:337–356
- [122] Heiser WJ (1995) Convergent computing by iterative majorization: theory and applications in multidimensional data analysis. In: Krzanowski WJ (ed) Recent advances in descriptive multivariate analysis. Clarendon, Oxford, pp 157–189
- [123] Henrici P (1982) Essentials of numerical analysis with pocket calculator demonstrations. Wiley, Hoboken
- [124] Herman GT (1980) Image reconstruction from projections: the fundamentals of computerized tomography. Springer, New York
- [125] Hestenes MR (1981) Optimization theory: the finite dimensional case. Robert E Krieger Publishing, Huntington
- [126] Hestenes MR, Karush WE (1951) A method of gradients for the calculation of the characteristic roots and vectors of a real symmetric matrix. *J Res Natl Bur Stand* 47:471–478
- [127] Hestenes MR, Stiefel E (1952) Methods of conjugate gradients for solving linear systems. *J Res Natl Bur Stand* 29:409–439
- [128] Higham NJ (2008) Functions of matrices: theory and computation. SIAM, Philadelphia
- [129] Hille E (1959) Analytic function theory, vol 1. Blaisdell, New York
- [130] Hiriart-Urruty J-B (1986) When is a point x satisfying $\nabla f(x) = \mathbf{0}$ a global minimum of $f(x)$? *Am Math Mon* 93:556–558
- [131] Hiriart-Urruty J-B, Claude Lemaréchal C (2001) Fundamentals of convex analysis. Springer, New York

- [132] Hochstadt H (1986) *The functions of mathematical physics*. Dover, New York
- [133] Hoel PG, Port SC, Stone CJ (1971) *Introduction to probability theory*. Houghton Mifflin, Boston
- [134] Hoffman K (1975) *Analysis in Euclidean space*. Prentice-Hall, Englewood Cliffs
- [135] Hoffman K, Kunze R (1971) *Linear algebra*, 2nd edn. Prentice-Hall, Englewood Cliffs
- [136] Horn RA, Johnson CR (1985) *Matrix analysis*. Cambridge University Press, Cambridge
- [137] Horn RA, Johnson CR (1991) *Topics in matrix analysis*. Cambridge University Press, Cambridge
- [138] Householder AS (1975) *The theory of matrices in numerical analysis*. Dover, New York
- [139] Hrusa W, Troutman JL (1981) Elementary characterization of classical minima. *Am Math Mon* 88:321–327
- [140] Hunter DR (2004) MM algorithms for generalized Bradley-Terry models. *Ann Stat* 32:386–408
- [141] Hunter DR, Lange K (2000) Quantile regression via an MM algorithm. *J Comput Graph Stat* 9:60–77
- [142] Hunter DR, Lange K (2004) A tutorial on MM algorithms. *Am Stat* 58:30–37
- [143] Hunter DR, Li R (2005) Variable selection using MM algorithms. *Ann Stat* 33:1617–1642
- [144] Jamshidian M, Jennrich RI (1995) Acceleration of the EM algorithm by using quasi-Newton methods. *J Roy Stat Soc B* 59:569–587
- [145] Jamshidian M, Jennrich RI (1997) Quasi-Newton acceleration of the EM algorithm. *J Roy Stat Soc B* 59:569–587
- [146] Jennrich RI, Moore RH (1975) Maximum likelihood estimation by means of nonlinear least squares. In: *Proceedings of the statistical computing section*. American Statistical Association, Washington, DC, pp 57–65
- [147] Jia R-Q, Zhao H, Zhao W (2009) Convergence analysis of the Bregman method for the variational model of image denoising. *Appl Comput Harmon Anal* 27:367–379

- [148] Karlin S, Taylor HM (1975) A first course in stochastic processes, 2nd edn. Academic, New York
- [149] Karush W (1939) Minima of functions of several variables with inequalities as side conditions. Master's Thesis, Department of Mathematics, University of Chicago, Chicago
- [150] Keener JP (1993) The Perron-Frobenius theorem and the ranking of football teams. *SIAM Rev* 35:80–93
- [151] Kelley CT (1999) Iterative methods for optimization. SIAM, Philadelphia
- [152] Khalfan HF, Byrd RH, Schnabel RB (1993) A theoretical and experimental study of the symmetric rank-one update. *SIAM J Optim* 3:1–24
- [153] Kiers HAL (1997) Weighted least squares fitting using ordinary least squares algorithms. *Psychometrika* 62:251–266
- [154] Kingman JFC (1993) Poisson processes. Oxford University Press, Oxford
- [155] Komiya H (1988) Elementary proof for Sion's minimax theorem. *Kodai Math J* 11:5–7
- [156] Kosowsky JJ, Yuille AL (1994) The invisible hand algorithm: solving the assignment problem with statistical physics. *Neural Network* 7:477–490
- [157] Kruskal JB (1964) Nonmetric multidimensional scaling: a numerical method. *Psychometrika* 29:115–129
- [158] Kruskal JB (1965) Analysis of factorial experiments by estimating monotone transformations of the data. *J Roy Stat Soc B* 27:251–263
- [159] Ku HH, Kullback S (1974) Log-linear models in contingency table analysis. *Biometrics* 10:452–458
- [160] Kuhn S (1991) The derivative á la Carathéodory. *Am Math Mon* 98:40–44
- [161] Kuhn HW, Tucker AW (1951) Nonlinear programming. In: Proceedings of the second Berkeley symposium on mathematical statistics and probability. University of California Press, Berkeley
- [162] Lange K (1994) An adaptive barrier method for convex programming. *Meth Appl Anal* 1:392–402

- [163] Lange K (1995) A gradient algorithm locally equivalent to the EM algorithm. *J Roy Stat Soc B* 57:425–437
- [164] Lange K (1995) A quasi-Newton acceleration of the EM algorithm. *Stat Sin* 5:1–18
- [165] Lange K (2002) *Mathematical and statistical methods for genetic analysis*, 2nd edn. Springer, New York
- [166] Lange K (2010) *Numerical analysis for statisticians*, 2nd edn. Springer, New York
- [167] Lange K, Carson R (1984) EM reconstruction algorithms for emission and transmission tomography. *J Comput Assist Tomogr* 8:306–316
- [168] Lange K, Fessler JA (1995) Globally convergent algorithms for maximum a posteriori transmission tomography. *IEEE Trans Image Process* 4:1430–1438
- [169] Lange K, Wu T (2008) An MM algorithm for multicategory vertex discriminant analysis. *J Comput Graph Stat* 17:527–544
- [170] Lange K, Zhou H (2012) MM algorithms for geometric and sigmoidal programming. *Math Program, Series A*, DOI 10.1007/s10107-012-0612-1
- [171] Lange K, Hunter D, Yang I (2000) Optimization transfer using surrogate objective functions (with discussion). *J Comput Graph Stat* 9:1–59
- [172] Lax PD (2007) *Linear algebra and its applications*, 2nd edn. Wiley, Hoboken
- [173] Ledoita O, Wolf M (2004) A well-conditioned estimator for large-dimensional covariance matrices. *J Multivar Anal* 88:365–411
- [174] Lee DD, Seung HS (1999) Learning the parts of objects by non-negative matrix factorization. *Nature* 401:788–791
- [175] Lee DD, Seung HS (2001) Algorithms for non-negative matrix factorization. *Adv Neural Inf Process Syst* 13:556–562
- [176] Lehmann EL (1986) *Testing statistical hypotheses*, 2nd edn. Wiley, Hoboken
- [177] Levina E, Rothman A, Zhu J (2008) Sparse estimation of large covariance matrices via a nested lasso penalty. *Ann Appl Stat* 2:245–263
- [178] Li Y, Arce GR (2004) A maximum likelihood approach to least absolute deviation regression. *EURASIP J Appl Signal Process* 2004:1762–1769

- [179] Little RJA, Rubin DB (1987) *Statistical analysis with missing data*. Wiley, Hoboken
- [180] Louis TA (1982) Finding the observed information matrix when using the EM algorithm. *J Roy Stat Soc B* 44:226–233
- [181] Luce RD (1959) *Individual choice behavior: a theoretical analysis*. Wiley, Hoboken
- [182] Luce RD (1977) The choice axiom after twenty years. *J Math Psychol* 15:215–233
- [183] Luenberger DG (1984) *Linear and nonlinear programming*, 2nd edn. Addison-Wesley, Reading
- [184] Magnus JR, Neudecker H (1988) *Matrix differential calculus with applications in statistics and econometrics*. Wiley, Hoboken
- [185] Maher MJ (1982) Modelling association football scores. *Stat Neerl* 36:109–118
- [186] Mangasarian OL, Fromovitz S (1967) The Fritz John necessary optimality conditions in the presence of equality and inequality constraints. *J Math Anal Appl* 17:37–47
- [187] Mardia KV, Kent JT, Bibby JM (1979) *Multivariate analysis*. Academic, New York
- [188] Marsden JE, Hoffman MJ (1993) *Elementary classical analysis*, 2nd edn. W H Freeman & Co, New York
- [189] Mazumder R, Hastie T, Tibshirani R (2010) Spectral regularization algorithms for learning large incomplete matrices. *J Mach Learn Res* 11:2287–2322
- [190] McLachlan GJ, Do K-A, Ambroise C (2004) *Analyzing microarray gene expression data*. Wiley, Hoboken
- [191] McLachlan GJ, Krishnan T (2008) *The EM algorithm and extensions*, 2nd edn. Wiley, Hoboken
- [192] McLachlan GJ, Peel D (2000) *Finite mixture models*. Wiley, Hoboken
- [193] McLeod RM (1980) *The generalized Riemann integral*. Mathematical Association of America, Washington, DC
- [194] McShane EJ (1973) The Lagrange multiplier rule. *Am Math Mon* 80:922–925

- [195] Meyer RR (1976) Sufficient conditions for the convergence of monotonic mathematical programming algorithms. *J Comput Syst Sci* 12:108–121
- [196] Michelot C (1986) A finite algorithm for finding the projection of a point onto the canonical simplex in \mathbb{R}^n . *J Optim Theor Appl* 50:195–200
- [197] Miller KS (1987) Some eclectic matrix theory. Robert E Krieger Publishing, Malabar
- [198] Moré JJ, Sorensen DC (1983) Computing a trust region step. *SIAM J Sci Stat Comput* 4:553–572
- [199] Narayanan A (1991) Algorithm AS 266: maximum likelihood estimation of the parameters of the Dirichlet distribution. *Appl Stat* 40:365–374
- [200] Nazareth L (1979) A relationship between the BFGS and conjugate gradient algorithms and its implications for new algorithms. *SIAM J Numer Anal* 16:794–800
- [201] Nedelman J, Wallenius T (1986) Bernoulli trials, Poisson trials, surprising variances, and Jensen’s inequality. *Am Stat* 40:286–289
- [202] Nelder JA, Wedderburn RWM (1972) Generalized linear models. *J Roy Stat Soc A* 135:370–384
- [203] Nemirovski AS, Todd MJ (2008) Interior-point methods for optimization. *Acta Numerica* 17:191–234
- [204] Nocedal J (1991) Theory of algorithms for unconstrained optimization. *Acta Numerica* 1991:199–242
- [205] Nocedal J, Wright S (2006) Numerical optimization, 2nd edn. Springer, New York
- [206] Orchard T, Woodbury MA (1972) A missing information principle: theory and applications. In: Proceedings of the 6th Berkeley symposium on mathematical statistics and probability. University of California Press, Berkeley, pp 697–715
- [207] Ortega JM (1990) Numerical analysis: a second course. Society for Industrial and Applied Mathematics, Philadelphia
- [208] Osher S, Burger M, Goldfarb D, Xu J, Yin W (2005) An iterative regularization method for total variation based image restoration. *Multiscale Model Simul* 4:460–489

- [209] Osher S, Mao T, Dong B, Yin W (2011) Fast linearized Bregman iteration for compressive sensing and sparse denoising. *Comm Math Sci* 8:93–111
- [210] Park MY, Hastie T (2008) Penalized logistic regression for detecting gene interactions. *Biostatistics* 9:30–50
- [211] Pauca VP, Piper J, Plemmons RJ (2006) Nonnegative matrix factorization for spectral data analysis. *Linear Algebra Appl* 416:29–47
- [212] Peressini AL, Sullivan FE, Uhl JJ Jr (1988) *The mathematics of nonlinear programming*. Springer, New York
- [213] Polya G (1954) *Induction and analogy in mathematics. Volume I of mathematics and plausible reasoning*. Princeton University Press, Princeton
- [214] Portnoy S, Koenker R (1997) The Gaussian hare and the Laplacian tortoise: computability of squared-error versus absolute-error estimators. *Stat Sci* 12:279–300
- [215] Press WH, Teukolsky SA, Vetterling WT, Flannery BP (1992) *Numerical recipes in Fortran: the art of scientific computing*, 2nd edn. Cambridge University Press, Cambridge
- [216] Rabiner L (1989) A tutorial on hidden Markov models and selected applications in speech recognition. *Proc IEEE* 77:257–285
- [217] Ranola JM, Ahn S, Sehl ME, Smith DJ, Lange K (2010) A Poisson model for random multigraphs. *Bioinformatics* 26:2004–2011
- [218] Rao CR (1973) *Linear statistical inference and its applications*, 2nd edn. Wiley, Hoboken
- [219] Robertson T, Wright FT, Dykstra RL (1988) *Order restricted statistical inference*. Wiley, Hoboken
- [220] Romano G (1995) New results in subdifferential calculus with applications to convex analysis. *Appl Math Optim* 32:213–234
- [221] Rockafellar RT (1996) *Convex analysis*. Princeton University Press, Princeton
- [222] Royden HL (1988) *Real analysis*, 3rd edn. Macmillan, London
- [223] Rudin W (1979) *Principles of mathematical analysis*, 3rd edn. McGraw-Hill, New York
- [224] Rudin LI, Osher S, Fatemi E (1992) Nonlinear total variation based noise removal algorithms. *Physica D* 60:259–268

- [225] Rustagi JS (1976) Variational methods in statistics. Academic, New York
- [226] Ruszczyński A (2006) Nonlinear optimization. Princeton University Press, Princeton
- [227] Sabatti C, Lange K (2002) Genomewide motif identification using a dictionary model. *Proc IEEE* 90:1803–1810
- [228] Sagan H (1969) Introduction to the calculus of variations. McGraw-Hill, New York
- [229] Santosa F, Symes WW (1986) Linear inversion of band-limited reflection seismograms. *SIAM J Sci Stat Comput* 7:1307–1330
- [230] Schmidt M, van den Berg E, Friedlander MP, Murphy K (2009) Optimizing costly functions with simple constraints: a limited-memory projected quasi-Newton algorithm. In: van Dyk D, Welling M (eds) Proceedings of The twelfth international conference on artificial intelligence and statistics (AISTATS), vol 5, pp 456–463
- [231] Schölkopf B, Smola AJ (2002) Learning with kernels: support vector machines, regularization, optimization, and beyond. MIT, Cambridge
- [232] Seber GAF, Lee AJ (2003) Linear regression analysis, 2nd edn. Wiley, Hoboken
- [233] Segel LA (1977) Mathematics applied to continuum mechanics. Macmillan, New York
- [234] Seneta E (1973) Non-negative matrices: an introduction to theory and applications. Wiley, Hoboken
- [235] Sha F, Saul LK, Lee DD (2003) Multiplicative updates for nonnegative quadratic programming in support vector machines. In: Becker S, Thrun S, Obermayer K (eds) Advances in neural information processing systems 15. MIT, Cambridge, pp 1065–1073
- [236] Silvapulle MJ, Sen PK (2005) Constrained statistical inference. Wiley, Hoboken
- [237] Sinkhorn R (1967) Diagonal equivalence to matrices with prescribed row and column sums. *Am Math Mon* 74:402–405
- [238] Sion M (1958) On general minimax theorems. *Pac J Math* 8:171–176
- [239] Smith CAB (1957) Counting methods in genetical statistics. *Ann Hum Genet* 21:254–276

- [240] Smith DR (1974) Variational methods in optimization. Dover, Mineola
- [241] Sorensen DC (1997) Minimization of a large-scale quadratic function subject to spherical constraints. *SIAM J Optim* 7:141–161
- [242] Srebro N, Jaakkola T (2003) Weighted low-rank approximations. In: Machine learning international workshop conference 2003. AAAI Press, 20:720–727
- [243] Steele JM (2004) The Cauchy-Schwarz master class: an introduction to the art of inequalities. Cambridge University Press and the Mathematical Association of America, Cambridge
- [244] Stein EM, Shakarchi R (2003) Complex analysis. Princeton University Press, Princeton
- [245] Stern RJ, Wolkowicz H (1995) Indefinite trust region subproblems and nonsymmetric eigenvalue perturbations. *SIAM J Optim* 5:286–313
- [246] Stoer J, Bulirsch R (2002) Introduction to numerical analysis, 3rd edn. Springer, New York
- [247] Strang G (1986) Introduction to applied mathematics. Wellesley-Cambridge, Wellesley
- [248] Strang G (1986) The fundamental theorem of linear algebra. *Am Math Mon* 100:848–855
- [249] Strang G (2003) Introduction to linear algebra, 3rd edn. Wellesley-Cambridge, Wellesley
- [250] Swartz C, Thomson BS (1988) More on the fundamental theorem of calculus. *Am Math Mon* 95:644–648
- [251] Tanner MA (1993) Tools for statistical inference: methods for the exploration of posterior distributions and likelihood functions, 2nd edn. Springer, New York
- [252] Taylor H, Banks SC, McCoy JF (1979) Deconvolution with the ℓ_1 norm. *Geophysics* 44:39–52
- [253] Teboulle M (1992) Entropic proximal mappings with applications to nonlinear programming. *Math Oper Res* 17:670–690
- [254] Theobald CM (1975) An inequality for the trace of the product of two symmetric matrices. *Math Proc Camb Phil Soc* 77:265–267

- [255] Thompson HB (1989) Taylor's theorem using the generalized Riemann integral. *Am Math Mon* 96:346–350
- [256] Tibshirani R (1996) Regression shrinkage and selection via the lasso. *J Roy Stat Soc B* 58:267–288
- [257] Tibshirani R, Saunders M, Rosset S, Zhu J, Knight K (2005) Sparsity and smoothness via the fused lasso. *J Roy Stat Soc B* 67:91–108
- [258] Tikhomirov VM (1990) Stories about maxima and minima. American Mathematical Society, Providence
- [259] Titterton DM, Smith AFM, Makov UE (1985) Statistical analysis of finite mixture distributions. Wiley, Hoboken
- [260] Trefethen LN, Bau D (1997) Numerical linear algebra. SIAM, Philadelphia
- [261] Uherka DJ, Sergott AM (1977) On the continuous dependence of the roots of a polynomial on its coefficients. *Am Math Mon* 84:368–370
- [262] Vandenberghe L, Boyd S, Wu S (1998) Determinant maximization with linear matrix inequality constraints. *SIAM J Matrix Anal Appl* 19:499–533
- [263] Van Ruitenburg J (2005) Algorithms for parameter estimation in the Rasch model. Measurement and Research Department Reports 2005–4. CITO, Arnhem
- [264] Vapnik V (1995) The nature of statistical learning theory. Springer, New York
- [265] Vardi Y, Shepp LA, Kaufman L (1985) A statistical model for positron emission tomography. *J Am Stat Assoc* 80:8–37
- [266] Von Neumann J (1928) Zur theorie der gesellschaftsspiele. *Math Ann* 100:295–320
- [267] Wang L, Gordon MD, Zhu J (2006) Regularized least absolute deviations regression and an efficient algorithm for parameter tuning. In: Proceedings of the sixth international conference on data mining (ICDM'06). IEEE Computer Society, Washington, DC, pp 690–700
- [268] Wang S, Yehya N, Schadt EE, Wang H, Drake TA, Lusk AJ (2006) Genetic and genomic analysis of a fat mass trait with complex inheritance reveals marked sex specificity. *PLoS Genet* 2:148–159

- [269] Watson GA (1992) Characterization of the subdifferential of some matrix norms. *Linear Algebra Appl* 170:1039–1053
- [270] Weeks DE, Lange K (1989) Trials, tribulations, and triumphs of the EM algorithm in pedigree analysis. *IMA J Math Appl Med Biol* 6:209–232
- [271] Weiszfeld E (1937) On the point for which the sum of the distances to n given points is minimum. *Ann Oper Res* 167:741 (Translated from the French original [Tohoku Math J 43:335–386 (1937)] and annotated by Frank Plastria)
- [272] Weston J, Elisseeff A, Schölkopf B, Tipping M (2003) Use of the zero-norm with linear models and kernel methods. *J Mach Learn Res* 3:1439–1461
- [273] Whyte BM, Gold J, Dobson AJ, Cooper DA (1987) Epidemiology of acquired immunodeficiency syndrome in Australia. *Med J Aust* 147:65–69
- [274] Wright MH (2005) The interior-point revolution in optimization: history, recent developments, and lasting consequences. *Bull Am Math Soc* 42:39–56
- [275] Wu CF (1983) On the convergence properties of the EM algorithm. *Ann Stat* 11:95–103
- [276] Wu TT, Lange K (2008) Coordinate descent algorithms for lasso penalized regression. *Ann Appl Stat* 2:224–244
- [277] Wu TT, Lange K (2010) Multicategory vertex discriminant analysis for high-dimensional data. *Ann Appl Stat* 4:1698–1721
- [278] Yee PL, Vyborný R (2000) *The integral: an easy approach after Kurzweil and Henstock*. Cambridge University Press, Cambridge
- [279] Yin W, Osher S, Goldfarb D, Darbon J (2008) Bregman iterative algorithms for ℓ_1 -minimization with applications to compressed sensing. *SIAM J Imag Sci* 1:143–168
- [280] Zhang Z, Lange K, Ophoff R, Sabatti C (2010) Reconstructing DNA copy number by penalized estimation and imputation. *Ann Appl Stat* 4:1749–1773
- [281] Zhou H, Lange K (2009) On the bumpy road to the dominant mode. *Scand J Stat* 37:612–631

- [282] Zhou H, Lange K (2010) MM algorithms for some discrete multivariate distributions. *J Comput Graph Stat* 19:645–665
- [283] Zhou H, Lange K (2012) A path algorithm for constrained estimation. *J Comput Graph Stat* DOI 10.1080/10618600.2012.681248
- [284] Zhou H, Lange K (2012) Path following in the exact penalty method of convex programming (submitted)

Index

- ABO genetic locus, 189
- Active constraint, 107
- Adaptive barrier methods,
 - 318–325
 - linear programming, 320
 - logarithmic, 318–320
- Admixtures, *see* EM algorithm,
 - cluster analysis
- Affine function, 108
 - convexity, 143
 - Fenchel conjugate, 345
 - majorization, 208
- Allele frequency estimation,
 - 189–191, 225–226
- Alternating least squares, 181
- Analytic function, 87
- ANOVA, 217
- Apollonius’s problem, 130
- Arithmetic–geometric mean
 - inequality, 145, 189
- Arithmetic–geometric mean
 - inequality, 2–3, 8
- Armijo rule, 303
- Attenuation coefficient, 199
- Backward algorithm, Baum’s,
 - 235
- Ball, 31
- Barrier method, 314–317
 - adaptive, 318–325
- Basis pursuit, 416, 433, 440
- Baum’s algorithms, 234–236
- Bayesian prior, 325, 371
- Bernstein polynomial, 169
- Binomial distribution, 159
 - score and information, 257
- Birkhoff’s theorem, 481
- Bivariate normal distribution
 - missing data, 239
- Block relaxation, 171–183
 - Bradley–Terry model, 210
 - canonical correlations, 175
 - global convergence of,
 - 302–303

- Block relaxation (*cont.*)
 - iterative proportional fitting, 177
 - k-means clustering, 174
 - local convergence, 296–297, 307
 - Maher’s sports model, 173
 - Sinkhorn’s algorithm, 172
- Blood type data, 256
- Blood type genes, 189, 225, 237
- Boundary point, 32
- Bounded set, 30
- Bradley–Terry model, 193, 210
- Bregman distance, 318
- Bregman iteration, 432–440
 - linearized, 435
 - split, 437
- Broyden–Fletcher–Goldfarb–Shanno update, 282

- Canonical correlations, 175–176, 182
- Carathéodory’s theorem, 141
- Cauchy sequence, 27
- Cauchy–Schwarz inequality, 7–8, 159
 - generalized, 348
- Censored variable, 239
- Chain rule, 85
 - for convex functions, 360
 - for second differential, 123
- Characteristic polynomial, 39
- Chebyshev’s inequality, 159
- Chernoff bound, 169
- Cholesky decomposition, 167, 476
- Closed function, 342
- Closed set, 30
- Closure, 33
- Cluster analysis, 174–175, 226–228, 325, 334
- Coercive function, 297–299, 310
- Coloring, 198
- Compact set, 33
 - subdifferential, 354
- Complete, 27
- Completeness, 27
 - and existence of suprema, 28
- Concave function, 9, 143
 - constraint qualification, 155
- Cone
 - finitely generated, 31, 477–479
 - normal, 356
 - polar, 350, 477–479
 - polyhedral, 477–479
 - problems, 378
 - tangent, 356
- Conjugate gradient algorithm, 275–278
- Conjugate vectors, 275
- Connected set, 44
 - arcwise, 44
- Constrained optimization, 339
- Contingency table
 - three-way, 177
- Continuous function, 34
- Continuously differentiable function, 88
- Contraction map, 389–393
- Convergence of optimization algorithms
 - local, 297
- Convergent sequence, 26
- Convex cone, 31
 - polar, 350
- Convex function, 9, 142
 - continuity, 149
 - differences of, 147, 491
 - differentiation rules, 358–364
 - directional derivatives, 150
 - integral, 152
 - minimization of, 152–158

- Convex hull, 141, 163, 342, 359, 363, 478, 497
- Convex programming, 318–325
 - convergence of MM
 - algorithm, 321–325
 - dual programs, *see* dual programs
 - Dijkstra’s algorithm, 384–388
 - for a geometric program, 320
 - linear classification, 403–406
- Convex regression, 386
- Convex set, 138–142
- Coordinate descent, 328–334
- Coronary disease data, 183
- Cousin’s lemma, 54
- Covariance matrix, 325, 370
- Critical point, 3
- Cross-validation, 328
- Cyclic coordinate descent
 - local convergence, 296–297
 - minimum of a quadratic, 182
 - saddle point, convergence to, 180
- Danskin’s theorem, 492
- Davidon’s formula, 281
- Davidon-Fletcher-Powell update, 283
- De Pierro majorization, 188
- Derivative
 - directional, 80
 - elementary functions, 76
 - equality of mixed partials, 80–81
 - forward directional, 80
 - partial, 79, 117–123
 - univariate, 56
- Descent direction, 249
- Differential, 81–88, 93–98
 - Carathéodory’s definition, 82–83
 - Fréchet’s definition, 82
- Gâteaux, 102
 - higher order, 117–123
 - linear function, 84
 - matrix-valued functions, 93–98
 - multilinear map, 84
 - quadratic function, 84
 - rules for calculating, 84, 85, 94–98, 122, 123
 - semidifferential, 487–497
 - subdifferential, 341
- Directional derivative, 80
 - rules for calculating, 490
- Dirichlet distribution
 - score and information, 265
- Dirichlet-multinomial
 - distribution, 215
- Discriminant analysis, 333
- Distance to a set, 35, 144
 - semidifferential, 495
 - subdifferential, 356
- Doubly stochastic matrix, 480
- Dual norm, 348
- Dual programs, 393–402
 - Duffin’s counterexample, 401
 - Fenchel conjugate, 395
 - geometric programming, 398
 - linear programming, 396
 - quadratic programming, 396–397
 - semidefinite programming, 397
- Duality gap, 395
- Duodenal ulcer blood type data, 256
- Dijkstra’s algorithm, 384–388, 399, 410
- Edgeworth’s algorithm, 329–330, 335
- Eigenvalue, 13
 - algorithm, 289
 - block relaxation, 297

- Eigenvalue (*cont.*)
 condition number, 289
 continuity, 39
 convexity, 146
 cyclic coordinate descent,
 182
 dominant, 25
 minimum, 126
 Rayleigh quotient, 182, 364
 subdifferential, 364
- Eigenvector, 13
 algorithm, 289
- Elsner–Koltracht–Neumann
 theorem, 390
- EM algorithm, 221–244
 allele frequency estimation,
 225
 ascent property, 222–224
 bivariate normal
 parameters, 239
 cluster analysis, 226–228,
 325–334
 E step, 222
 estimating binomial
 parameter, 240
 estimating multinomial
 parameters, 241
 exponential family, 237
 factor analysis, 231–234
 linear regression with right
 censoring, 239
 local convergence
 sublinear rate, 308
 M step, 222
 movie rating, 243–244
 sublinear convergence rate,
 307
 transmission tomography,
 228–230
 zero truncated data,
 242–243
- Entropy, 236
 Fenchel conjugate, 344
- Epigraph, 144, 342
- Equality constraint, 107
- Essential domain, 342
- Euclidean norm, 23–24
- Exact penalty method, 421–432
 convex programming, 424
- Expected information, 254
 admixture density, 266
 exponential families, 257
- Exponential distribution
 score and information, 257
- Exponential family, 224–225,
 255–256
 EM algorithm, 237
 expected information, 256
 generalized linear models,
 258
 natural, 432
- Extremal value, 3
 distinguishing from a saddle
 point, 124
- Extreme point, 481, 497
- Factor analysis, 231
- Factor loading matrix, 232
- Fan’s inequality, 365, 371, 483
- Farkas’ lemma, 130, 140–141
- Feasible point, 107
- Fenchel biconjugate, 345
- Fenchel conjugate, 6–7, 342–351
 ℓ_p function, 6
 affine function, 345
 differentiability, 495
 entropy, 344
 indicator, 347
 log determinant, 351
 problems, 375–378
 quadratic, 344
 subdifferential, 355
- Fenchel–Moreau theorem, 345
- Fenchel–Young inequality, 344,
 345
- Fermat’s principle, 9, 83, 341,
 355
- Fixed point, 292, 389
- Fletcher–Reeves update, 277
- Forward algorithm, Baum’s, 235

- Forward directional derivative,
 - 80
 - and subdifferentials,
 - 351–352
 - as a support function, 354
 - at a minimum, 153
 - of a maximum, 86
 - rules for calculating, 490
 - semidifferential, 487
 - sublinearity, 352
- Free variable, 108
- Frobenius matrix norm, 24
- Frobenius norm, 497
- Function
 - affine, 108
 - closed, 342
 - coercive, 297–299, 310
 - concave, 9, 143
 - continuous, 34
 - continuously differentiable,
 - 88
 - convex, 9, 142
 - differentiable, 81
 - distance to a set, 35
 - Gamma, 148
 - homogeneous, 101, 348
 - Huber’s, 267
 - indicator, 347
 - Lagrangian, 11, 373
 - link, 258
 - Lipschitz, 102, 149, 489
 - log posterior, 201
 - log-convex, 147
 - loglikelihood, 12, 156, 237
 - majorizing, 186
 - matrix exponential, 29–30
 - objective, 107
 - perspective, 346
 - potential, 201
 - rational, 35
 - Riemann’s zeta, 163
 - Rosenbrock’s, 19, 205
 - slope, 82
 - sublinear, 348
 - sublinear, 352
 - support, 347–348
 - uniformly continuous, 39
- Fundamental theorem of algebra, 38
- Fundamental theorem of calculus, 62–64
- Fundamental theorem of linear algebra, 356
- Fundamental theorem of linear programming, 401
- Gamma distribution
 - maximum likelihood estimation, 271
- Gamma function, 148
- Gauge function, 54
- Gauge integral, 53–54
- Gauss-Lucas theorem, 163
- Gauss-Newton algorithm, 252
 - as scoring, 257–258
- Gene expression, 331
- Generalized linear model,
 - 258–259, 332
- Geometric programming,
 - 157–158, 320
 - dual, 398
- Gibbs prior, 201
- Gibbs’ lemma, 132, 385
- Golden search, 279
- Gosper’s formula, 271
- Gradient algorithms, 303–306
- Gradient direction, 10
- Gradient vector, 9
- Hadamard semidifferential,
 - 487–497
- Hadamard’s inequality, 134
- Halfspace, 31
- Hardy, Littlewood, and Pólya inequality, 482
- Hardy-Weinberg law, 189
- Heine-Borel Theorem, 55
- Hermite interpolation, 278
- Hessian matrix, 9
- Hestenes-Stiefel update, 277

- Hidden Markov chain
 - EM algorithm, 235
- Hidden trials
 - binomial, 240
 - EM algorithm for, 241
 - multinomial, 241
 - Poisson or exponential, 241
- Higher-order differentials, 117–123
- Holder's inequality, 132, 162, 348
- Homogeneous function, 101, 348
- Huber's function, 267
- Hyperplane, 12, 31
- Image denoising, 437
- Implicit function theorem, 91–92
- Inactive constraint, 107
- Indicator function, 347
- Induced matrix norm, 25
- Inequality
 - arithmetic-geometric mean, 2–3, 8, 145
 - Cauchy-Schwarz, 7–8, 159
 - Chebyshev's, 159
 - Fan's, 365, 371, 483
 - Fenchel-Young, 345
 - Hölder's, 132, 162, 348
 - Hadamard's, 134
 - Hardy, Littlewood, and Pólya, 482
 - information, 223
 - Jensen's, 160
 - Lipschitz, 150
 - Markov's, 159
 - Minkowski's triangle, 170
 - Schlömilch's, 161–162
 - Young's, 376
- Inequality constraint, 107
- Infimal convolution, 377
- Information inequality, 223
- Integration by parts, 65
- Interior, 32
- Intermediate value theorem, 45, 55
- Inverse function theorem, 89–91
- Isotone regression, 386, 389, 430
- Iterative proportional fitting, 177–178, 183
- Jensen's inequality, 160
- K-means clustering, 174–175, 209
- K-medians clustering, 183
- Karush-Kuhn-Tucker theory
 - Kuhn-Tucker constraint qualification, 114–115
 - multiplier rule, *see* Lagrange multiplier rule
 - sufficient condition for a minimum, 125–128
- Kepler's problem, 4
- Kullback–Leibler distance, 376
- Kullback–Leibler divergence, 432
- Kullback-Leibler distance, 165
- L'Hôpital's rule, 99, 135
- Lagrange multiplier rule, 109–111, 372–375
- Lagrangian function, 11, 373, 393–395
- Lasso, 327–334
- Least absolute deviation regression, 192–193, 209, 329–330
- Least squares, 9–10, 179, 386
 - alternating, 181
 - isotone, 430
 - nonlinear, 252–253
 - nonnegative, 182, 210
 - right-censored data, 239
 - weighted, 218
- Least squares estimation, 330
- Leibnitz's formula, 99
- Limit inferior, 28
- Limit superior, 28
- Line search methods, 278–280
- Linear classification, 403–406
- Linear convergence, 292

- Linear logistic regression, 194, 211
- Linear programming, 108, 112, 320
 - dual for, 396
 - fundamental theorem, 401
- Linear regression, 179, 315
- Link function, 259
- Lipschitz inequality, 150
- Log posterior function, 201
- Log-convex function, 147
- Logarithmic barrier method, 318–320
- Loglikelihood function, 12, 156, 237
- Loglinear model, 177
 - observed information, 183
- Lower semicontinuity, 42–44, 342
- Luce’s ranking model, 219
- Maehly’s algorithm, 263
- Maher’s sport model, 173–174
- Majorizing function, 186, 204–206
- Mangasarian-Fromovitz
 - constraint qualification, 108, 116
- Markov chain
 - hidden, 234–236
 - stationary distribution, 141, 392
- Markov’s inequality, 159
- Marquardt’s method, 269
- Matrix
 - continuity considerations, 26
 - covariance estimation, 325, 370
 - eigenvalues of a symmetric, 13
 - exponential function, 29
 - factor loading, 232
 - Hessian, 9
 - induced norm, 25
 - nilpotent, 49
 - observed information, 13
 - positive definite, 36
 - rank semicontinuity, 44
 - skew-symmetric, 49
 - square root, 264
- Matrix completion problem, 368
- Matrix exponential function, 29–30
 - and differential equations, 77
- Matrix logarithm, 78
- Maximum likelihood estimation
 - allele frequency, 189
 - Dirichlet distribution, 251–252
 - exponential distribution, 253–254
 - hidden Markov chains, *see* Markov chain
 - multinomial distribution, 12–13, 235–325
 - multivariate normal
 - distribution, 156
 - Poisson distribution, 253
 - power series family, for a, 268
- Maxwell-Boltzmann distribution, 237
- Mean value theorem, 361
 - failure of, 89
 - multivariate, 88
 - univariate, 77
- Median, 361
- Method of false position, 278
- Michélot’s algorithm, 409
- Minkowski’s triangle inequality, 170
- Minkowski–Weyl theorem, 478
- Minorizing function, 187, 204–206
- Missing data
 - EM algorithm, 222, 234, 235
- Mixtures, *see* EM algorithm, cluster analysis
- MM algorithm, 185–219
 - acceleration, 259–311

- MM algorithm (*cont.*)
 allele frequency estimation, 189
 ANOVA, 217
 Bradley-Terry model, 193
 convergence for convex program, 321–325
 descent property, 186
 Dirichlet-multinomial, 215
 for discriminant analysis, 334
 geometric programming, 194
 global convergence of, 299–302
 linear logistic regression, 194
 linear regression, 191–193
 Luce’s model, 219
 majorization, 187–189
 matrix completion, 368
 motif finding, 215
 movie rating, 243–244
 random multigraph, 202–203
 transmission tomography, *see* transmission tomography
 zero-truncated data, 242–243
- MM gradient algorithm, 250–252
 convergence of, 294–295
 convex programming, 319
 Dirichlet distribution(, 251
 Dirichlet distribution), 252
- Model selection, 327–334
- Moreau-Rockafellar theorem, 360
- Motif finding, 215
- Movie rating, 243–244
- Multilinear map, 41–42, 51
 differential, 84
 second differential, 120
- Multilogit model, 272
- Multinomial distribution, 12, 190, 315, 325
 score and information, 257
- Multivariate normal
 distribution, 475–476
 maximum entropy, 236
 maximum likelihood, 156
- Neighborhood, 32
- Newton’s method, 245–259
 convergence of, 293–294
 least squares estimation, 252–253
 MM gradient algorithm, *see* MM gradient algorithm
 quadratic function, for, 265
 random multigraph, 249
 root finding, 246–248
 scoring, *see* scoring
 transmission tomography, 251
- Nilpotent matrix, 49
- Nonnegative matrix
 factorization, 179–180, 337–338
- Norm, 23–26
 ℓ_1 , 24, 490
 ℓ_∞ , 24, 492
 $\ell_{1,2}$, 381, 409
 dual, 348
 equivalence of, 37
 Euclidean, 23–24
 Frobenius matrix, 24
 induced matrix, 25, 494
 nuclear, 350, 381, 494
 orthogonally invariant, 493
 semidifferentiability, 489–494
 spectral, 25, 494
 subdifferential, 356
 trace, 350
- Normal cone, 356
- Normal distribution, 473–476
 mixtures, 226, 325
 multivariate, 475–476
 univariate, 473–474
- Normal equation, 9
- Nuclear norm, 350, 381

- Objective function, 107
- Observed information, 245
- Observed information matrix, 13
- Obtuse angle criterion, 140, 153
- Open set, 32
- Optimal experimental design, 413
- Order statistics, 364
- Orthogonal projection, 209

- Paracontractive map, 389–393
- Partial derivative, 79
- Penalized estimation, 327–334
- Penalty method, 314–317
- Perron-Frobenius theorem, 167
- Perspective of a function, 346
- Pixel, 200
- Poisson admixture model, 238
- Poisson distribution
 - contingency table data, modeling, 177
 - score and information, 257
- Poisson process, 197
- Polak–Ribière update, 277
- Polar cone, 350, 477–479
- Polyhedral set, 360, 477–480
- Polytope, 497
- Pool adjacent violators, 388–389
- Population genetics, *see* allele
 - frequency estimation inference of
 - maternal/paternal alleles in offspring, 14–16
- Positron emission tomography, 214
- Posterior mode, 201
- Posynomial, 158, 194–196, 211–213
- Potential function, 201
- Power plant problem, 334, 403
- Power series family, 268
- Primal program
 - convex, 395
- Projected gradient algorithm, 416–421
- Projection operators, 384, 410
- Proposition
 - Birkhoff, 481
 - Bolzano-Weierstrass, 33
 - Carathéodory, 141
 - Cousin, 55
 - Danskin, 492
 - Du Bois-Reymond, 457
 - Ekeland, 115
 - Fenchel-Moreau, 345
 - Gordon, 116, 148
 - Hahn-Banach, 353
 - Heine, 40
 - Jensen, 160
 - Liapunov, 301
 - Minkowski–Weyl, 478
 - Moreau-Rockafellar, 360
 - Ostrowski, 292
 - von Neumann–Fan, 483
 - Weierstrass, 37, 55
- q quantile, 207, 361
- QR decomposition, 475
- Quadratic bound principle, 188
- Quadratic convergence, 292
- Quadratic programming, 113, 114, 386
 - dual for, 396–397
- Quantal response model, 266
- Quantile regression, 362
- Quasi-convexity, 158
- Quasi-Newton algorithms, 281–284
 - ill-conditioning, avoiding, 288–289
- Random multigraph, 202–203, 249
- Random thinning, 198
- Rate of convergence, 292
- Rayleigh quotient, 182, 364
- Recurrence relations
 - hidden Markov chain, 235

- Relative topology, 34
- Riemann's zeta function, 163
- Rigid motion, 40–41
- Robust regression, 192, 266, 311
- Roots of a polynomial, 38
- Rosenbrock's function, 19, 205

- Saddle point, 3, 6, 124, 135, 302, 394
- Schlömilch's inequality, 161–162
- Score, 245
 - admixture density, 266
 - exponential families, 257
- Scoring, 254–257, 259
 - allele frequency estimation, 256
 - convergence of, 296
 - Gauss-Newton algorithm, 257–258
- Secant condition, 281
 - inverse, 283
- Second differential, 9, 127
 - chain rule for, 123
 - in optimization, 123
- Semicontinuity, 42–44
- Semidifferential, 487–497
- Set
 - arcwise connected, 44
 - closed, 30
 - compact, 33
 - connected, 44
 - convex, 138
 - indicator, 347
 - normal cone, 356
 - open, 32
 - polar cone, 350
 - polyhedral, 360, 477–480
 - tangent cone, 356
- Shadow values, 112
- Sherman-Morrison formula, 249, 282
 - Woodbury's generalization, 288
- Shrinkage, 362
- Signomial programming, 196–197, 213

- Simultaneous projection, 391
- Singular value decomposition, 181
- Sinkhorn's algorithm, 172–173, 181
- Skew-symmetric matrix, 49
- Slack variable, 108
- Slater constraint qualification, 155, 373
- Slope function, 82
- Snell's law, 4, 205
- Spectral functions, 365–372
- Spectral radius, 25
- Sphere, 31
- Stationary point, 3, 322
- Steepest ascent, 254
- Steepest descent, 10, 328
- Stirling's formula, 271
- Stopping criteria, 280
- Subdifferential, 341, 351–375
 - distance to a set, 356
 - eigenvalue, 364
 - Fenchel conjugate, 355
 - Fermat's principle, 355
 - Hahn-Banach theorem, 353
 - indicator function, 356
 - Lagrange multiplier rule, 375
 - mean value theorem, 361
 - median, 361
 - norm, 356
 - nuclear norm, 381
 - order statistics, 364
 - problems, 379–380
 - rules for forming, 358–364
- Subgradient, 151, 341
- Sublinear function, 348, 352
- Support function, 347–348
 - directional derivative, 354
- Survival analysis, 269–271
- Sylvester's criterion, 134

- Tangent cone, 356
- Tangent vectors and curves, 93
- Taylor expansion, 65
 - multivariate, 117–123

- Trace norm, [350](#)
- Transmission tomography,
[198–202](#), [228–230](#), [240](#)
- Trust region, [285](#), [413](#)

- Uniform convergence, [46](#)
- Uniformly continuous function,
[40](#)
- Univariate normal distribution,
see normal distribution,
univariate
- Upper semicontinuity, [42–44](#)

- Variance matrix, [325](#), [370](#)
- Von Neumann’s minimax
theorem, [402](#)

- Weierstrass approximation
theorem, [160](#)
- Weierstrass M-test, [47](#)
- Weiszfeld’s algorithm, [209](#)
- Woodbury’s formula, [288](#)

- Young’s inequality, [376](#)

- Zero-truncated data, [242–243](#)