

Endre Süli and David Mayers

An Introduction to Numerical Analysis

CAMBRIDGE

CAMBRIDGE

more information - www.cambridge.org/9780521810265

This page intentionally left blank

An Introduction to Numerical Analysis

Endre Süli and David F. Mayers

University of Oxford



CAMBRIDGE
UNIVERSITY PRESS

CAMBRIDGE UNIVERSITY PRESS

Cambridge, New York, Melbourne, Madrid, Cape Town, Singapore, São Paulo

Cambridge University Press

The Edinburgh Building, Cambridge CB2 2RU, United Kingdom

Published in the United States of America by Cambridge University Press, New York

www.cambridge.org

Information on this title: www.cambridge.org/9780521810265

© Cambridge University Press, 2003

This book is in copyright. Subject to statutory exception and to the provision of relevant collective licensing agreements, no reproduction of any part may take place without the written permission of Cambridge University Press.

First published in print format 2003

ISBN-13 978-0-511-07653-4 eBook (EBL)

ISBN-10 0-511-07653-3 eBook (EBL)

ISBN-13 978-0-521-81026-5 hardback

ISBN-10 0-521-81026-4 hardback

ISBN-13 978-0-521-00794-8 paperback

ISBN-10 0-521-00794-1 paperback

Cambridge University Press has no responsibility for the persistence or accuracy of URLs for external or third-party internet websites referred to in this book, and does not guarantee that any content on such websites is, or will remain, accurate or appropriate.

Contents

<i>Preface</i>	<i>page</i>	vii
1	Solution of equations by iteration	1
1.1	Introduction	1
1.2	Simple iteration	2
1.3	Iterative solution of equations	17
1.4	Relaxation and Newton's method	19
1.5	The secant method	25
1.6	The bisection method	28
1.7	Global behaviour	29
1.8	Notes	32
	Exercises	35
2	Solution of systems of linear equations	39
2.1	Introduction	39
2.2	Gaussian elimination	44
2.3	LU factorisation	48
2.4	Pivoting	52
2.5	Solution of systems of equations	55
2.6	Computational work	56
2.7	Norms and condition numbers	58
2.8	Hilbert matrix	72
2.9	Least squares method	74
2.10	Notes	79
	Exercises	82
3	Special matrices	87
3.1	Introduction	87
3.2	Symmetric positive definite matrices	87
3.3	Tridiagonal and band matrices	93

3.4	Monotone matrices	98
3.5	Notes	101
	Exercises	102
4	Simultaneous nonlinear equations	104
4.1	Introduction	104
4.2	Simultaneous iteration	106
4.3	Relaxation and Newton's method	116
4.4	Global convergence	123
4.5	Notes	124
	Exercises	126
5	Eigenvalues and eigenvectors of a symmetric matrix	133
5.1	Introduction	133
5.2	The characteristic polynomial	137
5.3	Jacobi's method	137
5.4	The Gerschgorin theorems	145
5.5	Householder's method	150
5.6	Eigenvalues of a tridiagonal matrix	156
5.7	The QR algorithm	162
5.7.1	The QR factorisation revisited	162
5.7.2	The definition of the QR algorithm	164
5.8	Inverse iteration for the eigenvectors	166
5.9	The Rayleigh quotient	170
5.10	Perturbation analysis	172
5.11	Notes	174
	Exercises	175
6	Polynomial interpolation	179
6.1	Introduction	179
6.2	Lagrange interpolation	180
6.3	Convergence	185
6.4	Hermite interpolation	187
6.5	Differentiation	191
6.6	Notes	194
	Exercises	195
7	Numerical integration – I	200
7.1	Introduction	200
7.2	Newton–Cotes formulae	201
7.3	Error estimates	204
7.4	The Runge phenomenon revisited	208
7.5	Composite formulae	209

7.6	The Euler–Maclaurin expansion	211
7.7	Extrapolation methods	215
7.8	Notes	219
	Exercises	220
8	Polynomial approximation in the ∞-norm	224
8.1	Introduction	224
8.2	Normed linear spaces	224
8.3	Best approximation in the ∞ -norm	228
8.4	Chebyshev polynomials	241
8.5	Interpolation	244
8.6	Notes	247
	Exercises	248
9	Approximation in the 2-norm	252
9.1	Introduction	252
9.2	Inner product spaces	253
9.3	Best approximation in the 2-norm	256
9.4	Orthogonal polynomials	259
9.5	Comparisons	270
9.6	Notes	272
	Exercises	273
10	Numerical integration – II	277
10.1	Introduction	277
10.2	Construction of Gauss quadrature rules	277
10.3	Direct construction	280
10.4	Error estimation for Gauss quadrature	282
10.5	Composite Gauss formulae	285
10.6	Radau and Lobatto quadrature	287
10.7	Note	288
	Exercises	288
11	Piecewise polynomial approximation	292
11.1	Introduction	292
11.2	Linear interpolating splines	293
11.3	Basis functions for the linear spline	297
11.4	Cubic splines	298
11.5	Hermite cubic splines	300
11.6	Basis functions for cubic splines	302
11.7	Notes	306
	Exercises	307

12	Initial value problems for ODEs	310
12.1	Introduction	310
12.2	One-step methods	317
12.3	Consistency and convergence	321
12.4	An implicit one-step method	324
12.5	Runge–Kutta methods	325
12.6	Linear multistep methods	329
12.7	Zero-stability	331
12.8	Consistency	337
12.9	Dahlquist’s theorems	340
12.10	Systems of equations	341
12.11	Stiff systems	343
12.12	Implicit Runge–Kutta methods	349
12.13	Notes	353
	Exercises	355
13	Boundary value problems for ODEs	361
13.1	Introduction	361
13.2	A model problem	361
13.3	Error analysis	364
13.4	Boundary conditions involving a derivative	367
13.5	The general self-adjoint problem	370
13.6	The Sturm–Liouville eigenvalue problem	373
13.7	The shooting method	375
13.8	Notes	380
	Exercises	381
14	The finite element method	385
14.1	Introduction: the model problem	385
14.2	Rayleigh–Ritz and Galerkin principles	388
14.3	Formulation of the finite element method	391
14.4	Error analysis of the finite element method	397
14.5	<i>A posteriori</i> error analysis by duality	403
14.6	Notes	412
	Exercises	414
	An overview of results from real analysis	419
	WWW-resources	423
	<i>Bibliography</i>	424
	<i>Index</i>	429

Preface

This book has grown out of printed notes which accompanied lectures given by ourselves and our colleagues over many years to undergraduate mathematicians at Oxford. During those years the contents and the arrangement of the lectures have changed substantially, and this book has a wider scope than is currently taught. It contains mathematics which, in an ideal world, would be part of the equipment of any well-educated mathematician.

Numerical analysis is the branch of mathematics concerned with the theoretical foundations of numerical algorithms for the solution of problems arising in scientific applications. The subject addresses a variety of questions ranging from the approximation of functions and integrals to the approximate solution of algebraic, transcendental, differential and integral equations, with particular emphasis on the stability, accuracy, efficiency and reliability of numerical algorithms. The purpose of this book is to provide an elementary introduction into this active and exciting field, and is aimed at students in the second year of a university mathematics course.

The book addresses a wide range of numerical problems in algebra and analysis. Chapter 2 deals with the solution of systems of linear equations, a process which can be completed in a finite number of arithmetical operations. In the rest of the book the solution of a problem is sought as the limit of an infinite sequence; in that sense the output of the numerical algorithm is an ‘approximate’ solution. This need not, however, mean any relaxation of the usual standards of rigorous analysis. The idea of convergence of a sequence of real numbers (x_n) to a real number x is very familiar: given any positive value of ϵ there exists a positive integer N such that $|x_n - x| < \epsilon$ for all n such that $n > N$. In such a situation one can obtain as accurate an approximation to x as

required by calculating sufficiently many members of the sequence, or just one member, sufficiently far along. A ‘pure mathematician’ would prefer the exact answer, π , but the sorts of guaranteed accurate approximations which will be discussed here are entirely satisfactory in real-life applications.

Numerical analysis brings two new ideas to the usual discussion of convergence of sequences. First, we need, not just the existence of N , but a good estimate of how large it is; and it may be too large for practical calculations. Second, rather than being asked for the limit of a given sequence, we are usually given the existence of the limit (or its approximate location on the real line) and then have to construct a sequence which converges to it. If the rate of convergence is slow, so that the value of N is large, we must then try to construct a better sequence, one that converges to π more rapidly. These ideas have direct applications in the solution of a single nonlinear equation in Chapter 1, the solution of systems of nonlinear equations in Chapter 4 and the calculation of the eigenvalues and eigenvectors of a matrix in Chapter 5.

The next six chapters are concerned with polynomial approximation, and show how, in various ways, we can construct a polynomial which approximates, as accurately as required, a given continuous function. These ideas have an obvious application in the evaluation of integrals, where we calculate the integral of the approximating polynomial instead of the integral of the given function.

Finally, Chapters 12 to 14 deal with the numerical solution of ordinary differential equations, with Chapter 14 presenting the fundamentals of the finite element method. The results of Chapter 14 can be readily extended to linear second-order partial differential equations.

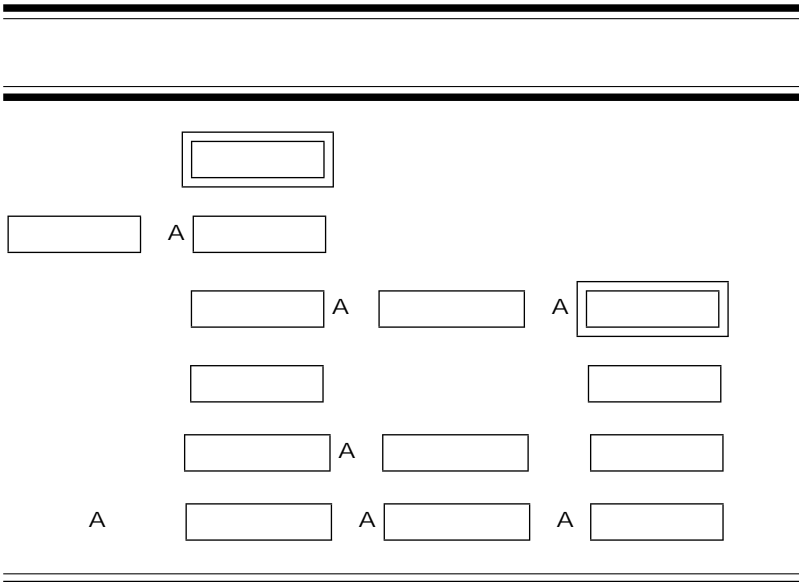
We have tried to make the coverage as complete as is consistent with remaining quite elementary. The limitations of size are most obvious in Chapter 12 on the solution of initial value problems for ordinary differential equations. This is an area where a number of excellent books are available, at least one of which is published in two weighty volumes. Chapter 12 does not describe or analyse anything approaching all the available methods, but we hope we have included some of those in most common use.

There is a selection of Exercises at the end of each chapter. All these exercises are theoretical; students are urged to apply all the methods described to some simple examples to see what happens. A few of the exercises will be found to require some heavy algebraic manipulation; these have been included because we assume that readers will have ac-

cess to some computer algebra system such as Maple or Mathematica, which then make the algebraic work almost trivial. Those involved in teaching courses based on this book may obtain copies of L^AT_EX files containing solutions to these exercises by applying to the publisher by email (). Although the material presented in this book does not presuppose the reader's acquaintance with mathematical software packages, the importance of these cannot be overemphasised. In Appendix B, a brief set of pointers is provided to relevant software repositories.

Our treatment is intended to maintain a reasonably high standard of rigour, with many theorems and formal proofs. The main prerequisite is therefore some familiarity with elementary real analysis. Appendix A lists the standard theorems (labelled **Theorem A.1**, **A.2**, . . . , **A7**) which are used in the book, together with proofs of one or two of them which might be less familiar. Some knowledge of basic matrix algebra is assumed. We have also used some elementary ideas from the theory of normed linear spaces in a number of places; complete definitions and examples are given. Some prior knowledge of these areas would be helpful, although not essential.

The chart below indicates how the chapters of the book are inter-related. They show, in particular, how Chapters 1 to 5 form a largely self-contained unit, as do Chapters 6 to 10.



We have included some historical notes throughout the book. As well as hoping to stimulate an interest in the development of the subject, these notes show how wide a historical range even this elementary book covers. Many of the methods were developed by the great mathematicians of the seventeenth and eighteenth centuries, including Newton, Euler and Gauss, but what is usually known as Gaussian elimination for the solution of systems of linear equations was known to the Chinese two thousand years ago. At the other end of the historical scale, the analysis of the eigenvalue problem, and the numerical solution of differential equations, are much more recent, and are due to mathematicians who are still very much alive. Many of our historical notes are based on the excellent biographical database at the history of mathematics website

We have tried to eradicate as many typographical errors from the text as possible; however, we are mindful that some may have escaped our attention. We plan to post any typos reported to us on

We wish to express our gratitude to Professor Bill Morton for setting us off on this *tour de force*, to David Tranah at Cambridge University Press for encouraging us to persist with the project, and to the staff of the Press for not only improving the appearance of the book and eliminating a number of typographical errors, but also for correcting and improving some of our mathematics. We also wish to thank our colleagues at the Oxford University Computing Laboratory, particularly Nick Trefethen, Mike Giles and Andy Wathen, for keeping our spirits up, and to Paul Houston at the Department of Mathematics and Computer Science of the University of Leicester for his help with the final example in the book.

Above all, we are grateful to our families for their patience, support and understanding: this book is dedicated to them.

Solution of equations by iteration

1.1 Introduction

Equations of various kinds arise in a range of physical applications and a substantial body of mathematical research is devoted to their study. Some equations are rather simple: in the early days of our mathematical education we all encountered the single *linear* equation $ax + b = 0$, where a and b are real numbers and $a \neq 0$, whose solution is given by the formula $x = -b/a$. Many equations, however, are *nonlinear*: a simple example is $ax^2 + bx + c = 0$, involving a quadratic polynomial with real coefficients a , b , c , and $a \neq 0$. The two solutions to this equation, labelled x_1 and x_2 , are found in terms of the coefficients of the polynomial from the familiar formulae

$$x_1 = \frac{-b + \sqrt{b^2 - 4ac}}{2a}, \quad x_2 = \frac{-b - \sqrt{b^2 - 4ac}}{2a}. \quad (1.1)$$

It is less likely that you have seen the more intricate formulae for the solution of cubic and quartic polynomial equations due to the sixteenth century Italian mathematicians Niccolo Fontana Tartaglia (1499–1557) and Lodovico Ferrari (1522–1565), respectively, which were published by Girolamo Cardano (1501–1576) in 1545 in his *Artis magna sive de regulis algebraicis liber unus*. In any case, if you have been led to believe that similar expressions involving radicals (roots of sums of products of coefficients) will supply the solution to any polynomial equation, then you should brace yourself for a surprise: no such closed formula exists for a general polynomial equation of degree n when $n \geq 5$. It transpires that for each $n \geq 5$ there exists a polynomial equation of degree n with

Theorem 1.1 *Let f be a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line. Assume, further, that $f(a)f(b) < 0$; then, there exists ξ in $[a, b]$ such that $f(\xi) = 0$.*

Proof If $f(a) = 0$ or $f(b) = 0$, then $\xi = a$ or $\xi = b$, respectively, and the proof is complete. Now, suppose that $f(a)f(b) < 0$. Then, $f(a)f(b) < 0$; in other words, 0 belongs to the open interval whose endpoints are $f(a)$ and $f(b)$. By the Intermediate Value Theorem (Theorem A.1), there exists ξ in the open interval (a, b) such that $f(\xi) = 0$. \square

To paraphrase Theorem 1.1, if a continuous function f has opposite signs at the endpoints of the interval $[a, b]$, then the equation $f(x) = 0$ has a solution in (a, b) . The converse statement is, of course, false. Consider, for example, a continuous function defined on $[a, b]$ which changes sign in the open interval (a, b) an even number of times, with $f(a)f(b) < 0$; then, $f(a)f(b) < 0$ even though $f(x) = 0$ has solutions inside $[a, b]$. Of course, in the latter case, there exist an even number of subintervals of (a, b) at the endpoints of each of which f *does* have opposite signs. However, finding such subintervals may not always be easy.

To illustrate this last point, consider the rather pathological function

$$f: x \mapsto \frac{1}{2} - \frac{1}{1 + Mx - 1.05}, \quad (1.2)$$

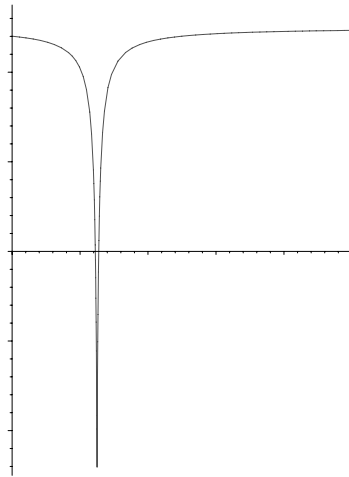
depicted in Figure 1.1 for x in the closed interval $[0.8, 1.8]$ and $M = 200$. The solutions $x = 1.05 - (1/M)$ and $x = 1.05 + (1/M)$ to the equation $f(x) = 0$ are only a distance $2/M$ apart and, for large and positive M , locating them computationally will be a challenging task.

Remark 1.1 *If you have access to the mathematical software package Maple, plot the function f by typing*

```
!"#$ % & ' ( ) * + , - . / : ; < = > ? @ [ \ ] ^ _ ` { | } ~
```

at the Maple command line, and then repeat this experiment by choosing $M = 2000, 20000, 200000, 2000000, \text{ and } 20000000$ in place of the number 200. What do you observe? For the last two values of M , replot the function f for x in the subinterval $[1.04999, 1.05001]$.

An alternative sufficient condition for the existence of a solution to the equation $f(x) = 0$ is arrived at by rewriting it in the equivalent form $x = g(x) = 0$ where g is a certain real-valued function, defined



$$\frac{1}{2} \frac{1}{1+200 \frac{1}{1.05}}$$

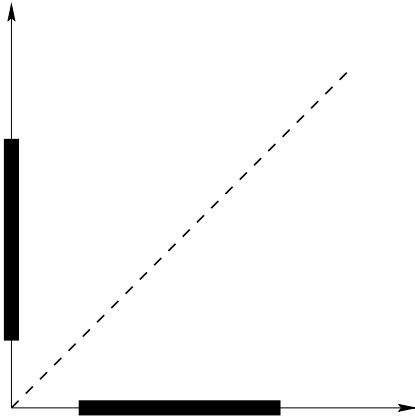
and continuous on $[a, b]$; the choice of g and its relationship with f will be clarified below through examples. Upon such a transformation the problem of solving the equation $f(x) = 0$ is converted into one of finding x such that $g(x) = x$.

Theorem 1.2 (Brouwer's Fixed Point Theorem) *Suppose that g is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and let $g(x) \in [a, b]$ for all $x \in [a, b]$. Then, there exists $x \in [a, b]$ such that $x = g(x)$; the real number x is called a **fixed point of the function g** .*

Proof Let $f(x) = x - g(x)$. Then, $f(a) = a - g(a) \geq 0$ since $g(a) \in [a, b]$ and $f(b) = b - g(b) \leq 0$ since $g(b) \in [a, b]$. Consequently, $f(a)f(b) \leq 0$, with f defined and continuous on the closed interval $[a, b]$. By Theorem 1.1 there exists $x \in [a, b]$ such that $0 = f(x) = x - g(x)$. \square

Figure 1.2 depicts the graph of a function $x = g(x)$, defined and continuous on a closed interval $[a, b]$ of the real line, such that $g(x)$ belongs to $[a, b]$ for all x in $[a, b]$. The function g has three fixed points in the interval $[a, b]$: the x -coordinates of the three points of intersection of the graph of g with the straight line $y = x$.

Of course, any equation of the form $f(x) = 0$ can be rewritten in the



Although the ability to verify the existence of a solution to the equation $f(x) = 0$ is important, none of what has been said so far provides a *method* for solving this equation. The following definition is a first step in this direction: it will lead to the construction of an algorithm for computing an approximation to the fixed point of the function g , and will thereby supply an approximate solution to the equivalent equation $f(x) = 0$.

Definition 1.1 *Suppose that g is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and assume that $g(x) \in [a, b]$ for all $x \in [a, b]$. Given that $x_0 \in [a, b]$, the recursion defined by*

$$x_k = g(x_{k-1}), \quad k = 0, 1, 2, \dots, \quad (1.3)$$

is called a **simple iteration**; the numbers $x_k, k \geq 0$, are referred to as **iterates**.

If the sequence (x_k) defined by (1.3) converges, the limit must be a fixed point of the function g , since g is continuous on a closed interval. Indeed, writing $x = \lim_{k \rightarrow \infty} x_k$, we have that

$$x = \lim_{k \rightarrow \infty} x_k = \lim_{k \rightarrow \infty} g(x_{k-1}) = g(\lim_{k \rightarrow \infty} x_{k-1}) = g(x), \quad (1.4)$$

where the second equality follows from (1.3) and the third equality is a consequence of the continuity of g .

A sufficient condition for the convergence of the sequence (x_k) is provided by our next result which represents a refinement of Brouwer's Fixed Point Theorem, under the additional assumption that the mapping g is a contraction.

Definition 1.2 (Contraction) *Suppose that g is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line. Then, g is said to be a **contraction** on $[a, b]$ if there exists a constant L such that $0 < L < 1$ and*

$$|g(x) - g(y)| \leq L|x - y| \quad x, y \in [a, b]. \quad (1.5)$$

Remark 1.2 *The terminology 'contraction' stems from the fact that when (1.5) holds with $0 < L < 1$, the distance $|g(x) - g(y)|$ between the images of the points x, y is (at least $1/L$ times) smaller than the distance*

x y between x and y . More generally, when L is any positive real number, (1.5) is referred to as a **Lipschitz condition**.

Armed with Definition 1.2, we are now ready to state the main result of this section.

Theorem 1.3 (Contraction Mapping Theorem) *Let g be a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line, and assume that $g(x) \in [a, b]$ for all $x \in [a, b]$. Suppose, further, that g is a contraction on $[a, b]$. Then, g has a unique fixed point in the interval $[a, b]$. Moreover, the sequence (x_n) defined by (1.3) converges to x^* as $k \rightarrow \infty$ for any starting value x_0 in $[a, b]$.*

Proof The existence of a fixed point for g is a consequence of Theorem 1.2. The uniqueness of this fixed point follows from (1.5) by contradiction: for suppose that g has a second fixed point, \tilde{x} , in $[a, b]$. Then,

$$|x^* - \tilde{x}| = |g(x^*) - g(\tilde{x})| \leq L|x^* - \tilde{x}|,$$

i.e., $(1 - L)|x^* - \tilde{x}| \leq 0$. As $1 - L > 0$, we deduce that $x^* = \tilde{x}$.

Let x_0 be any element of $[a, b]$ and consider the sequence (x_n) defined by (1.3). We shall prove that (x_n) converges to the fixed point x^* . According to (1.5) we have that

$$|x_n - x^*| = |g(x_{n-1}) - g(x^*)| \leq L|x_{n-1} - x^*|, \quad k \geq 1,$$

from which we then deduce by induction that

$$|x_n - x^*| \leq L^n |x_0 - x^*|, \quad k \geq 1. \tag{1.6}$$

As $L \in (0, 1)$, it follows that $\lim_{n \rightarrow \infty} L^n = 0$, and hence we conclude that $\lim_{n \rightarrow \infty} |x_n - x^*| = 0$. □

Let us illustrate the Contraction Mapping Theorem by an example.

Example 1.2 *Consider the equation $f(x) = 0$ on the interval $[1, 2]$ with $f(x) = e^{-2x} - 1$, as in Example 1.1. Recall from Example 1.1 that this equation has a solution, x^* , in the interval $[1, 2]$, and x^* is a fixed point of the function g defined on $[1, 2]$ by $g(x) = \ln(2x + 1)$.*

1 H : , + B# ' 6 #. * I B \$
 D C & H, # @ * F 6D ,
 6 6 1 1 , 9 6 , , , 6
 1 6 J G 6 , , ,
 6

Table 1.1. The sequence (x_k) defined by (1.8).

x_0	1
x_1	1.26
x_2	1.26
x_3	1.26
x_4	1.26
x_5	1.26
x_6	1.26
x_7	1.26
x_8	1.26
x_9	1.26
x_{10}	1.26

Now, the function g is defined and continuous on the interval $[1, 2]$, and g is differentiable on $(1, 2)$. Thus, by the Mean Value Theorem (Theorem A.3), for any x, y in $[1, 2]$ we have that

$$g(x) - g(y) = g'(c)(x - y) = g'(c) \cdot (x - y) \tag{1.7}$$

for some c that lies between x and y and is therefore in the interval $[1, 2]$. Further, $g'(x) = -2/(2x + 1)$ and $g''(x) = 4/(2x + 1)^2$. As $g'(x) < 0$ for all x in $[1, 2]$, g is monotonic decreasing on $[1, 2]$. Hence $g(1) > g(c) > g(2)$, *i.e.*, $g'(c) \in [2/5, 2/3]$. Thus we deduce from (1.7) that

$$|g(x) - g(y)| \leq L |x - y| \quad x, y \in [1, 2],$$

with $L = 2/3$. According to the Contraction Mapping Theorem, the sequence (x_k) defined by the simple iteration

$$x_{k+1} = \ln(2x_k + 1), \quad k = 0, 1, 2, \dots, \tag{1.8}$$

converges to α for any starting value x_0 in $[1, 2]$. Let us choose $x_0 = 1$, for example, and compute the next 11 iterates, say. The results are shown in Table 1.1. Even though we have carried six decimal digits, after 11 iterations only the first two decimal digits of the iterates x_k appear to have settled; thus it seems likely that $\alpha = 1.26$ to two decimal digits.

You may now wonder how many iterations we should perform in (1.8)

to ensure that all six decimals have converged to their correct values. In order to answer this question, we need to carry out some analysis.

Theorem 1.4 Consider the simple iteration (1.3) where the function g satisfies the hypotheses of the Contraction Mapping Theorem on the bounded closed interval $[a, b]$. Given $x_0 \in [a, b]$ and a certain tolerance $\epsilon > 0$, let $k(\epsilon)$ denote the smallest positive integer such that x_k is no more than ϵ away from the (unknown) fixed point α , i.e., $|x_k - \alpha| < \epsilon$, for all $k \geq k(\epsilon)$. Then,

$$k(\epsilon) = \frac{\ln(x_0 - \alpha) - \ln(\epsilon)}{\ln(1/L)} + 1, \tag{1.9}$$

where, for a real number x , $[x]$ signifies the largest integer less than or equal to x .

Proof From (1.6) in the proof of Theorem 1.3 we know that

$$x_k - \alpha = L(x_{k-1} - \alpha), \quad k \geq 1.$$

Using this result with $k = 1$, we obtain

$$\begin{aligned} x_1 - \alpha &= L(x_0 - \alpha) \\ x_2 - \alpha &= L(x_1 - \alpha) \\ &= L^2(x_0 - \alpha) \\ &\vdots \\ x_k - \alpha &= L^k(x_0 - \alpha). \end{aligned}$$

Hence

$$|x_k - \alpha| = \frac{1}{L^k} |x_0 - \alpha|.$$

By substituting this into (1.6) we get

$$|x_k - \alpha| = \frac{1}{L^k} |x_0 - \alpha| < \epsilon. \tag{1.10}$$

Thus, in particular, $|x_k - \alpha| < \epsilon$ provided that

$$L^k \frac{1}{L^k} |x_0 - \alpha| < \epsilon.$$

On taking the (natural) logarithm of each side in the last inequality, we find that $|x_k - \alpha| < \epsilon$ for all k such that

$$k > \frac{\ln(x_0 - \alpha) - \ln(\epsilon)}{\ln(1/L)}.$$

Therefore, the smallest integer $k(\epsilon)$ such that $|x_k - \alpha| < \epsilon$ for all

k $k(\epsilon)$ cannot exceed the expression on the right-hand side of the inequality (1.9). \square

This result provides an upper bound on the maximum number of iterations required to ensure that the error between the k th iterate x_k and the (unknown) fixed point α is below the prescribed tolerance ϵ . Note, in particular, from (1.9), that if L is close to 1, then $k(\epsilon)$ may be quite large for any fixed ϵ . We shall revisit this point later on in the chapter.

Example 1.3 Now we can return to Example 1.2 to answer the question posed there about the maximum number of iterations required, with starting value $x_0 = 1$, to ensure that the last iterate computed is correct to six decimal digits.

Letting $\epsilon = 0.5 \cdot 10^{-6}$ and recalling from Example 1.2 that $L = 2/3$, the formula (1.9) yields $k(\epsilon) = \lceil [32.778918] + 1 \rceil$, so we have that $k(\epsilon) = 33$. In fact, 33 is a somewhat pessimistic overestimate of the number of iterations required: computing the iterates x_k successively shows that already x_{30} is correct to six decimal digits, giving $\alpha = 1.256431$.

Condition (1.5) can be rewritten in the following equivalent form:

$$\left| \frac{g(x) - g(y)}{x - y} \right| \leq L \quad x, y \in [a, b], \quad x \neq y,$$

with $L \in (0, 1)$, which can, in turn, be rephrased by saying that the absolute value of the slope of the function g does not exceed L on $(0, 1)$. Assuming that g is a differentiable function on the open interval (a, b) , the Mean Value Theorem (Theorem A.3) tells us that

$$\frac{g(x) - g(y)}{x - y} = g'(\xi)$$

for some ξ that lies between x and y and is therefore contained in the interval (a, b) .

We shall therefore adopt the following assumption that is somewhat stronger than (1.5) but is easier to verify in practice:

$$g \text{ is differentiable on } (a, b) \text{ and} \tag{1.11}$$

$$L \in (0, 1) \text{ such that } |g'(x)| \leq L \text{ for all } x \in (a, b).$$

Consequently, Theorem 1.3 still holds when (1.5) is replaced by (1.11).

We note that the requirement in (1.11) that g be differentiable is

indeed more demanding than the Lipschitz condition (1.5): for example,
 $g(x) = x$

If the conditions of Theorem 1.5 are satisfied in the vicinity of a fixed point α , then the sequence (x_k) defined by the iteration $x_{k+1} = g(x_k)$, $k \geq 0$, will converge to α for any starting value x_0 that is sufficiently close to α . If, on the other hand, the conditions of Theorem 1.5 are violated, there is no guarantee that any sequence (x_k) defined by the iteration $x_{k+1} = g(x_k)$, $k \geq 0$, will converge to the fixed point α for any starting value x_0 near α . In order to distinguish between these two cases, we introduce the following definition.

Definition 1.3 *Suppose that g is a real-valued function, defined and continuous on the bounded closed interval $[a, b]$, such that $g(x) \in [a, b]$*

If (1.15) holds with $\mu = 0$, then the sequence (x_k) is said to converge to **superlinearly**.

If (1.15) holds with $\mu \in (0, 1)$ and $x_k = x^*$, $k = 0, 1, 2, \dots$, then (x_k) is said to converge to **linearly**, and the number $\mu = \log_{10} \mu$ is then called the **asymptotic rate of convergence** of the sequence. If (1.15) holds with $\mu = 1$ and $x_k = x^*$, $k = 0, 1, 2, \dots$, the rate of convergence is slower than linear and we say that the sequence converges to **sublinearly**.

The words ‘at least’ in this definition refer to the fact that we only have inequality in x_k , which may be all that can be ascertained in practice. Thus, it is really the sequence of bounds that converges linearly.

For a linearly convergent sequence the asymptotic rate of convergence measures the number of correct decimal digits gained in one iteration; in particular, the number of iterations required in order to gain one more correct decimal digit is at most $[1/\mu] + 1$. Here $[1/\mu]$ denotes the largest integer that is less than or equal to $1/\mu$.

Under the hypotheses of Theorem 1.5, the equalities (1.14) will hold with $\mu = g'(x^*) \in [0, 1)$, and therefore the sequence (x_k) generated by the simple iteration will converge to the fixed point x^* linearly or superlinearly.

Example 1.4 Given that α is a fixed positive real number, consider the function g defined on the interval $[0, 1]$ by

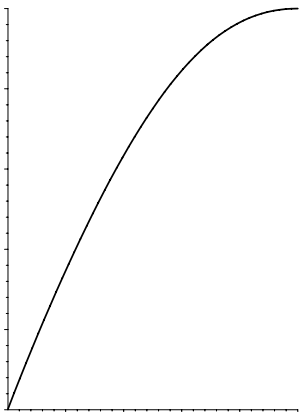
$$g(x) = \begin{cases} 2^{-x} & \text{for } 0 < x < 1, \\ 0 & \text{for } x = 0. \end{cases}$$

As $\lim_{x \rightarrow 0^+} g(x) = 0$, the function g is continuous on $[0, 1]$. Moreover, g is strictly monotonic increasing on $[0, 1]$ and $g(x) \in [0, 1/2] \subset [0, 1]$ for all x in $[0, 1]$. We note that $x = 0$ is a fixed point of g (cf. Figure 1.3).

Consider the sequence (x_k) defined by $x_{k+1} = g(x_k)$, $k \geq 0$, with $x_0 = 1$. It is a simple matter to show by induction that $x_k = 2^{-2^k}$, $k \geq 0$. Thus we deduce that (x_k) converges to $x^* = 0$ as $k \rightarrow \infty$. Since

$$\lim_{k \rightarrow \infty} \frac{x_{k+1}}{x_k} = \mu = \begin{cases} 1 & \text{for } 0 < x < 1, \\ - & \text{for } x = 1, \\ 0 & \text{for } x > 1, \end{cases}$$

we conclude that for $x_0 \in (0, 1)$ the sequence (x_k) converges to $x^* = 0$ sublinearly. For $x_0 = 1$ it converges to $x^* = 0$ linearly with asymptotic rate



of the sequence (x_k) is $\log |g'(x)|$. Evidently, a small value of $|g'(x)|$ corresponds to a large positive value of L and will result in more rapid convergence, while if $|g'(x)| < 1$ but $|g'(x)|$ is very close to 1, L will be a small positive number and the sequence will converge very slowly.

Next, we discuss the behaviour of the iteration (1.3) in the vicinity of an *unstable fixed point* α . If $|g'(\alpha)| > 1$, then the sequence (x_k) defined by (1.3) does not converge to α from any starting value x_0 ; the next theorem gives a rigorous proof of this fact.

Theorem 1.6 *Suppose that $\alpha = g(\alpha)$, where the function g has a continuous derivative in some neighbourhood of α , and let $|g'(\alpha)| > 1$. Then, the sequence (x_k) defined by $x_{k+1} = g(x_k)$, $k \geq 0$, does not converge to α from any starting value x_0 , $x_0 \neq \alpha$.*

Proof Suppose that $x_0 \neq \alpha$. As in the proof of Theorem 1.5, we can see that there is an interval $I = [\alpha - \delta, \alpha + \delta]$, $\delta > 0$, in which $|g'(x)| \leq L > 1$ for some constant L . If x_0 lies in this interval, then

$$|x_1 - \alpha| = |g(x_0) - g(\alpha)| = |x_0 - \alpha| |g'(\xi)| \leq L |x_0 - \alpha|,$$

for some ξ between x_0 and α . If x_1 lies in I the same argument shows that

$$|x_2 - \alpha| \leq L |x_1 - \alpha| \leq L^2 |x_0 - \alpha|,$$

and so on. Evidently, after a finite number of steps some member of the sequence x_0, x_1, x_2, \dots must be outside the interval I , since $L > 1$. Hence there can be no value of $k = k(\epsilon)$ such that $|x_k - \alpha| < \epsilon$ for all $k \geq k$, and the sequence therefore does not converge to α . \square

Example 1.5 *In this example we explore the simple iteration (1.3) for g defined by*

$$g(x) = -\frac{1}{2}(x + c)$$

where c is a fixed constant.

The fixed points of the function g are the solutions of the quadratic equation $x = -\frac{1}{2}(x + c)$, which are $1 - \frac{1}{2}c$. If $c > 1$ there are no solutions (in the set of real numbers, that is!), if $c = 1$ there is one solution in \mathbb{R} , and if $c < 1$ there are two.

1, 5, 4) #, 5, # 1, 5 \$
 B B D D 4 6, 5 G, 10 B 10 B D D

Suppose now that $c <$

Table 1.2. The sequences (u_n) and (v_n) in Example 1.6.

	\$	\$)
	\$)
	\$)	\$)
&	\$)	\$)
'	\$)	\$)
(\$)	\$)
%	\$)	\$)
)	\$)	\$)
\$	\$)	\$)
	\$)	\$)
	\$)	\$)

is still larger than 0.9. Although (u_n) eventually converges faster than v_n , we find that $u_n = (0.99)^n$ becomes smaller than $v_n = (k + 1)^{-n}$ when

$$k > \frac{10}{\ln(1/0.99)} \ln(k + 1).$$

This first happens when $k = 9067$, at which point u_n and v_n are both roughly 10^{-n} . In this rather extreme example the concept of asymptotic rate of convergence is not useful, since for any practical purposes (v_n) converges faster than (u_n) .

1.3 Iterative solution of equations

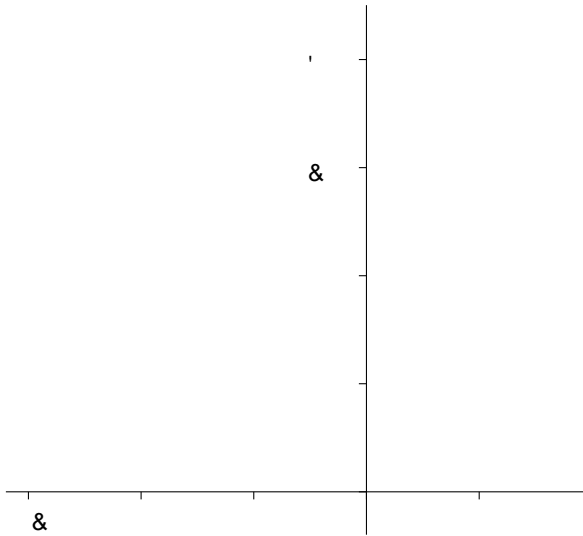
In this section we apply the idea of simple iteration to the solution of equations. Given a real-valued continuous function f , we wish to construct a sequence (x_n) , using iteration, which converges to a solution of $f(x) = 0$. We begin with an example where it is easy to derive various such sequences; in the next section we shall describe a more general approach.

Example 1.7 Consider the problem of determining the solutions of the equation $f(x) = 0$, where $f: x \mapsto e^x - x - 2$.

Since $f'(x) = e^x - 1$ the function f is monotonic increasing for positive x and monotonic decreasing for negative values of x . Moreover,

$$\begin{aligned}
 f(1) &= e^{-3} < 0, \\
 f(2) &= e^{-4} > 0, \\
 f(-1) &= e^{-1} < 0, \\
 f(-2) &= e^{-2} > 0.
 \end{aligned}
 \tag{1.16}$$

Hence the equation $f(x) = 0$ has exactly one positive solution, which lies in the interval $(1, 2)$, and exactly one negative solution, which lies in the interval $(-2, -1)$. This is illustrated in Figure 1.4, which shows the graphs of the functions $y = e^{-x}$ and $y = x + 2$ on the same axes. We shall write α for the positive solution and β for the negative solution.



will converge to the positive solution, α , provided that the starting value x_0 is sufficiently close to it. As $0 < g'(\alpha) < 1/3$, the asymptotic rate of convergence of (x_k) to α is certainly greater than $\log_3 3$.

On the other hand, $g'(\alpha) > 1$ since $2 < \alpha < 1$, so the sequence (x_k) defined by (1.17) cannot converge to the solution α . It is not difficult to prove that for $x_0 > \alpha$ the sequence (x_k) converges to α while if $x_0 < \alpha$ the sequence will decrease monotonically until $x_k = 2$ for some k , and then the iteration breaks down as $g(x_k)$ becomes undefined.

The equation $f(x) = 0$ may also be written in the form $x = e^{-2x}$, suggesting the sequence (x_k) defined by the iteration

$$x_{k+1} = e^{-2x_k}, \quad k = 0, 1, 2, \dots$$

In this case $g(x) = e^{-2x}$ and $g'(x) = -2e^{-2x}$. Hence $g'(\alpha) > 1$, $g'(\alpha) < e^{-2}$, showing that the sequence (x_k) may converge to α , but cannot converge to α . It is quite straightforward to show that the sequence converges to α for any $x_0 < \alpha$, but diverges to $+\infty$ when $x_0 > \alpha$.

As a third alternative, consider rewriting the equation $f(x) = 0$ as $x = g(x)$ where the function g is defined by $g(x) = x(e^{-x})/2$; the fixed points of the associated iteration $x_{k+1} = g(x_k)$ are the solutions α and β of $f(x) = 0$, and also the point 0. For this iteration neither of the fixed points, α or β , is stable, and the sequence (x_k) either converges to 0 or diverges to $+\infty$.

Evidently the given equation may be written in many different forms, leading to iterations with different properties.

1.4 Relaxation and Newton's method

In the previous section we saw how various ingenious devices lead to iterations which may or may not converge to the desired solutions of a given equation $f(x) = 0$. We would obviously benefit from a more generally applicable iterative method which would, except possibly in special cases, produce a sequence (x_k) that always converges to a required solution. One way of constructing such a sequence is by relaxation.

$$\begin{matrix} 1) & 1 & , & 6 & & 6 & & , & ' & & & & & & 6 \\ & & & 5 & < & = & & , & 1 & , & , & & & & G & , & B & D & E \\ & & & 6 & & & \text{B\#\#&\&D} & , & 5 & & 1 & 1 & , & 6 & 5 & 5 & & 04 \end{matrix}$$

Definition 1.5 Suppose that f is a real-valued function, defined and continuous in a neighbourhood of a real number α . **Relaxation** uses the sequence (x_k) defined by

$$x_{k+1} = x_k + \lambda (f(x_k) - \alpha), \quad k = 0, 1, 2, \dots, \quad (1.18)$$

where

corresponding to

$$= \frac{4}{2M + \dots}$$

On defining $g(x) = x - f(x)$, we then deduce that

$$g'(x) < 1, \quad x \in [a, b]. \tag{1.19}$$

Thus we can apply Theorem 1.5 to conclude that the sequence (x_k) defined by the relaxation iteration (1.18) converges to α , provided that x_0 is in the interval $[a, b]$. The asymptotic rate of convergence of the relaxation iteration (1.18) to α is at least $\log \dots$. \square

We can now extend the idea of relaxation by allowing f to be a continuous function of x in a neighbourhood of α rather than just a constant. This suggests an iteration

$$x_{k+1} = x_k - (x_k) f(x_k), \quad k = 0, 1, 2, \dots,$$

corresponding to a simple iteration with $g(x) = x - (x)f(x)$. If the sequence (x_k) converges, the limit α will be a solution of $f(x) = 0$, except possibly when $(\alpha) = 0$. Moreover, as we have seen, the ultimate rate of convergence is determined by $g'(\alpha)$. Since $f(\alpha) = 0$, it follows that $g'(\alpha) = 1 - (\alpha) f'(\alpha)$, and (1.19) suggest using a function ϕ which makes $1 - (\alpha) f'(\alpha)$ small. The obvious choice is $\phi(x) = 1/f(x)$, and leads us to Newton's method.

Definition 1.6 *Newton's method for the solution of $f(x) = 0$ is defined by*

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, 2, \dots, \tag{1.20}$$

with prescribed starting value x_0 . We implicitly assume in the defining formula (1.20) that $f'(x_k) \neq 0$ for all $k \geq 0$.

1) , *# ' , #& & : 4 % , , : , , 6 #? * , 6 #? , # ' , #& & 2 F , 5 6 5 , 1 , , 6 , 6 #> 4 , L , : +4 1 & 4) #& > 6 1 , , 1 6 \$, 6 B , 6 M , % L , E , 6 D , 6t , || , E , " 5 4 1

Newton's method is a simple iteration with $g(x) = x - f(x)/f'(x)$. Its geometric interpretation is illustrated in Figure 1.5: the tangent to the curve $y = f(x)$ at the point $(x, f(x))$ is the line with the equation $y = f(x)$

We note that unlike the definition of linear convergence where μ was required to belong to the interval $(0, 1)$, all we demand here is that $\mu > 0$. The reason is simple: when $q > 1$, (1.21) implies suitably rapid decay of the sequence (x_k) irrespective of the size of μ .

Example 1.8 Let $c > 1$ and $q > 1$. The sequence (x_k) defined by $x_{k+1} = c^{-k} x_k$, $k = 0, 1, 2, \dots$, converges to 0 with order q .

Theorem 1.8 (Convergence of Newton's method) Suppose that f is a continuous real-valued function with continuous second derivative f'' , defined on the closed interval $I = [a, b]$, $b - a > 0$, such that $f(a) = 0$ and $f'(a) \neq 0$. Suppose further that there exists a positive constant A such that

$$\left| \frac{f''(x)}{f'(x)} \right| \leq A \quad \forall x, y \in I.$$

If $x_0 \in I$, where h is the smaller of $b - a$ and $1/A$, then the sequence (x_k) defined by Newton's method (1.20) converges quadratically to α .

Proof Suppose that $x_0 \in I$, $h = \min\{b - a, 1/A\}$, so that $x_0 \in I$. Then, by Taylor's Theorem (Theorem A.4), expanding about the point $x_0 \in I$,

$$0 = f(\alpha) = f(x_0) + (x_0 - \alpha)f'(x_0) + \frac{(x_0 - \alpha)^2}{2}f''(\xi), \quad (1.22)$$

for some ξ between x_0 and α , and therefore in the interval I . Recalling (1.20), this shows that

$$x_1 - \alpha = -\frac{(x_0 - \alpha)^2 f''(\xi)}{2f'(x_0)}. \quad (1.23)$$

Since $x_0 \in I$, we have $x_1 \in I$. As we are given that $x_0 \in I$ it follows by induction that $x_k \in I$ for all $k \geq 0$; hence (x_k) converges to α as $k \rightarrow \infty$.

Now, ξ lies between x_0 and α , and therefore (ξ) also converges to α as $k \rightarrow \infty$. Since f' and f'' are continuous on I , it follows from (1.23) that

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \alpha}{x_k - \alpha} = \frac{f''(\alpha)}{2f'(\alpha)}, \quad (1.24)$$

which, according to Definition 1.7, implies quadratic convergence of the sequence (x_k) to α with $\mu = f''(\alpha)/2f'(\alpha)$, $\mu \in (0, A/2]$. \square

The conditions of the theorem implicitly require that $f'(x) \neq 0$, for otherwise the quantity $f'(x)/f'(y)$ could not be bounded in a neighbourhood of

1.5 The secant method

So far we have considered iterations which can be written in the form $x_{k+1} = g(x_k)$, $k \geq 0$, so that the new value is expressed in terms of the old one. It is also possible to define an iteration of the form $x_{k+1} = g(x_k, x_{k-1})$, $k \geq 1$, where the new value is expressed in terms of two previous values. In particular, we shall consider two applications of this idea, leading to the secant method and the method of bisection, respectively.

Remark 1.3 *We note in passing that one can consider more general iterative methods of the form*

$$x_{k+1} = g(x_k, x_{k-1}, \dots, x_{k-m}), \quad k = m, m+1, \dots,$$

with $m \geq 1$ fixed; here, we shall confine ourselves to the simplest case when $m = 1$ as this is already sufficiently illuminating.

Using Newton's method to solve a nonlinear equation $f(x) = 0$ requires explicit knowledge of the first derivative f' of the function f . Unfortunately, in many practical situations f' is not explicitly available or it can only be obtained at high computational cost. In such cases, the value $f'(x_k)$ in (1.20) can be approximated by a difference quotient; that is,

$$f'(x_k) \approx \frac{f(x_k) - f(x_{k-1})}{x_k - x_{k-1}}.$$

Replacing $f'(x_k)$ in (1.20) by this difference quotient leads us to the following definition.

Definition 1.8 *The secant method is defined by*

$$x_{k+1} = x_k - f(x_k) \frac{x_k - x_{k-1}}{f(x_k) - f(x_{k-1})}, \quad k = 1, 2, 3, \dots, \quad (1.25)$$

where x_0 and x_1 are given starting values. It is implicitly assumed here that $f(x_k) - f(x_{k-1}) \neq 0$ for all $k \geq 1$.

The method is illustrated in Figure 1.6. The new iterate x_{k+1} is obtained from x_{k-1} and x_k by drawing the chord joining the points $P(x_{k-1}, f(x_{k-1}))$ and $Q(x_k, f(x_k))$, and using as x_{k+1} the point at which this chord intersects the x -axis. If x_{k-1} and x_k are close together and f



Table 1.3. Comparison of the secant method and Newton's method for the solution of $e^x - 2 = 0$.

	1	/	0
	&		%&\$(&
&	&%%%%('%'
'	%" \$		'% \$&
((& \$\$		'% \$&
%	'()'(
)	'% \$		
	'% \$&		

where x_{k+1} is between x_k and x_{k-1} , and x_{k+2} lies between x_k and x_{k-1} . Hence, if $x_{k-1} > 1$ and $x_k > 1$, then also $x_{k+1} > 1$ and $x_{k+2} > 1$. Therefore,

$$x_{k+2} - 1 \leq \frac{5}{3} (x_k - 1) \quad (1.29)$$

Thus, $x_k > 1$ and the sequence (x_k) converges to α at least linearly, with rate at least $\log(3/2)$, provided that $x_0 > 1$ and $x_1 > 1$. \square

In fact, it can be shown that

$$\lim_{k \rightarrow \infty} \frac{x_{k+1} - \alpha}{x_k - \alpha} = \mu \quad (1.30)$$

where μ is a positive constant and $q = -(1 + \sqrt{5})/2 \approx -1.6$, so that the convergence of the sequence (x_k) to α is faster than linear, but not as fast as quadratic. (See Exercise 10.)

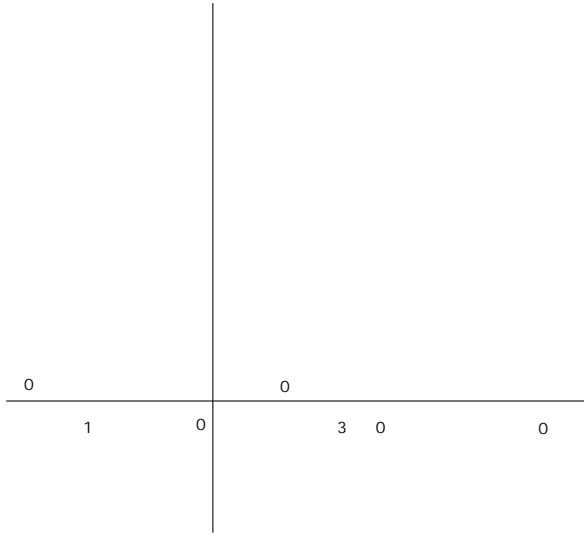
This is illustrated in Table 1.3, which compares two iterative methods for the solution of $f(x) = 0$ with $f: x \mapsto e^x - 2$; the first is the secant method, starting from $x_0 = 1$, $x_1 = 3$, while the second is Newton's method starting from $x_0 = 1$.

This experiment shows the faster convergence of Newton's method, but it must be remembered that each iteration of Newton's method requires the calculation of both $f(x_k)$ and $f'(x_k)$, while each iteration of the secant method requires the calculation of $f(x_k)$ only (as $f(x_{k-1})$ has already been computed). In our examples the computations are quite trivial, but in a practical situation the calculation of each value of $f(x_k)$ and $f'(x_k)$ may demand a substantial amount of work, and then

each iteration of Newton's method is likely to involve at least twice as much work as one iteration of the secant method.

1.6 The bisection method

Suppose that f is a real-valued function defined and continuous on a bounded closed interval $[a, b]$ of the real line and such that $f(a) < 0$ and $f(b) > 0$ for some $[a, b]$. A very simple iterative method for the solution of the nonlinear equation $f(x) = 0$ can be constructed by beginning with an interval $[a$



,) 2 , 0 0 ! 0 0 ! 0 0

$(b - a)/2$. The bisection method is therefore very robust, though Newton's method will always win once the current iterate is sufficiently close to α .

If the initial interval $[a, b]$ contains more than one solution, the limit of the bisection method will depend on the positions of these solutions. Figure 1.7 illustrates a possible situation, where $[a, b]$ contains three solutions. Since $f(c)$ has the same sign as $f(b)$ the second interval is $[a, c]$, and the sequence (c_n) of midpoints defined by (1.31) converges to the solution α_1 . If however the initial interval is $[a, b]$ the sequence of midpoints converges to the solution α_2 .

1.7 Global behaviour

We have already seen how an iteration will often converge to a limit if the starting value is sufficiently close to that limit. The behaviour of the iteration, when started from an arbitrary starting value, can be very complicated. In this section we shall consider two examples. No theorems will be stated: our aim is simply to illustrate various kinds of behaviour.

First consider the simple iteration defined by

$$x_{k+1} = g(x_k), \quad k = 0, 1, 2, \dots, \quad \text{where } g(x) = ax(1-x), \quad (1.33)$$

which is often known as the **logistic equation**. We require the constant a to lie in the range $0 < a \leq 4$, for then if the starting value x_0 is in the interval $[0, 1]$, then all members of the sequence (x_k) also lie in $[0, 1]$. The function g has two fixed points: $x = 0$ and $x = 1 - 1/a$. The fixed point at 0 is stable if $0 < a < 1$, and the fixed point at $1 - 1/a$ is stable if $1 < a < 3$. The behaviour of the iteration for these values of a is what might be expected from this information, but for larger values of the parameter a the behaviour of the sequence (x_k) becomes increasingly complicated.

For example, when $a = 3.4$ there is no stable fixed point, and from any starting point the sequence eventually oscillates between two values, which are 0.45 and 0.84 to two decimal digits. These are the two stable fixed points of the double iteration

$$x_{k+2} = g(g(x_k)), \quad g(g(x)) = g(g(x)) = a x(1-x)[1 - ax(1-x)]. \quad (1.34)$$

When $3 < a < 1 + \sqrt{6}$, the fixed points of g^2 are the two fixed points of g , that is 0 and $1 - 1/a$, and also

$$\frac{1}{2} \left(1 + \frac{1}{a} \pm \frac{1}{a} \sqrt{a^2 - 2a - 3} \right). \quad (1.35)$$

This behaviour is known as a stable two-cycle (see Exercise 12).

When $a > 1 + \sqrt{6}$ all the fixed points of g^2 are unstable. For example, when $a = 3.5$ all sequences (x_k) defined by (1.33) tend to a stable 4-cycle, taking successive values 0.50, 0.87, 0.38 and 0.83.

For larger values of the parameter a the sequences become chaotic. For example, when $a = 3.99$ there are no stable fixed points or limit-cycles, and the members of any sequence appear random. In fact it can be shown that for such values of a the members of the sequence are *dense* in a subinterval of $[0, 1]$: there exist real numbers α and β , $\alpha < \beta$, such that any subinterval of (α, β) , however small, contains an infinite subsequence of (x_k) . For the value $a = 3.99$ the maximal interval (α, β) is (0.00995, 0.99750) to five decimal digits. Starting from $x_0 = 0.75$ we find that the interval (0.70, 0.71), for example, contains the subsequence

$$x_1, x_2, x_3, x_4, x_5, \dots \quad (1.36)$$

The sequence does not show any apparent regular behaviour. The calculation is extremely sensitive: if we replace x_k by $x_k + \epsilon$, and write

that from any starting value Newton's method eventually converges to a solution, which might be α . However, it is certainly *not* true that the sequence converges to the solution closest to the starting point; indeed, if this were true, no sequence could converge to α . It is easy to see why the behaviour is much more complicated than this.

The Newton iteration converges to the solution at 0 from any point in the interval $(-0.327, 0.445)$. As we see from Figure 1.8, the iteration will converge exactly to 0 in one iteration if we start from the x -coordinate of any of the points a_1 , a_2 and a_3 ; at each of these three points the tangent to the curve passes through the origin. Since f is continuous, this means that there is an open interval surrounding each of these points from which the Newton iteration will converge to 0. The maximal such intervals are $(-1.555, -1.487)$, $(1.735, 1.817)$ and $(3.514, 3.529)$ to three decimal digits. In the same way, there are several points at which the tangent to the curve passes through the point $(A, 0)$, where A is the x -coordinate of the point a . Starting from one of these points, the Newton iteration will evidently converge exactly to the solution at 0 in two steps; surrounding each of these points there is an open interval from which the iteration will converge to 0.

Now suppose we define the sets S_m , $m = 1, 0, 1, 3, \dots$, where S_m consists of those points from which the Newton iteration converges to the zero at m . Then, an extension of the above argument shows that each of the sets S_m is the union of an infinite number of disjoint open intervals. The remarkable property of these sets is that, if α is a boundary point of one of the sets S_m , then it is also a boundary point of all the other sets as well. This means that any neighbourhood of such a point α , however small, contains an infinite number of members of each of the sets S_m . For example, we have seen that the iteration starting from any point in the interval $(-0.327, 0.445)$ converges to 0. We find that the end of this interval lies between 0.4457855 and 0.4457860; Table 1.4 shows the limits of various Newton iterations starting from points near this boundary. Each of these points is, of course, itself surrounded by an open interval which gives the same limit.

1.8 Notes

Theorem 1.2 is a special case of Brouwer's Fixed Point Theorem. Luitzen Egbertus Jan Brouwer (1881–1966) was professor of set theory, function theory and axiomatics at the University of Amsterdam, and made major contributions to topology. Brouwer was a mathematical genius with

Table 1.4. *Limit of Newton's method near a boundary point.*

0	3
"()'	
"()('	
"() (
"() ((
"() %	
"() %(
"())	
"())(
"())	
"() (
"() \$	
"() \$(&
"() \$	

strong mystical and philosophical leanings. For an historical overview of Brouwer's life and work we refer to the recent book of Dirk Van Dalen, *Mystic, Geometer, and Intuitionist. The Life of L.E.J. Brouwer: the Dawning Revolution*, Clarendon Press, Oxford, 1999.

The Contraction Mapping Theorem, as stated here, is a simplified version of Banach's fixed point theorem. Stefan Banach founded modern functional analysis and made outstanding contributions to the theory of topological vector spaces, measure theory, integration, the theory of sets, and orthogonal series. For an inspiring account of Banach's life and times, see R. Kaluza, *Through the Eyes of a Reporter: the Life of Stefan Banach*, Birkhäuser, Boston, MA, 1996.

In our definitions of linear convergence and convergence with order q , we followed Definitions 2.1 and 2.2 in Chapter 4 of

, *Numerical Analysis: an Introduction*, Birkhäuser, Boston, MA, 1997.

Exciting surveys of the history of Newton's method are available in T. Ypma, Historical development of the Newton-Raphson method, *SIAM Rev.* **37**, 531-551, 1995, H. Goldstine, *History of Numerical Analysis from the Sixteenth through the Nineteenth Century*, Springer, New York, 1977; and in Chapter 6 of Jean-Luc Chabert (Editor), *A History of Algorithms from the Pebble to the Microchip*, Springer, New York, 1999. As

¹ * ' , #. @ " - % CO 6 B D C *# % #@ >
: 5 5 " B D4

is noted in these sources, Newton's *De analysi per aequationes numero terminorum infinitas*, probably dating from mid-1669, is sometimes regarded as the historical source of the method, despite the fact that, surprisingly, there is no trace in this tract of the familiar recurrence relation $x_{n+1} = x_n - f(x_n)/f'(x_n)$ bearing Newton's name, nor is there a mention of the idea of derivative. Instead, the paper contains an example of a cubic polynomial whose roots are found by purely algebraic and rather complicated substitutions. In 1690, Joseph Raphson (1648–1715) in the Preface to his *Analysis aequationum universalis* describes his version of Newton's method as 'not only, I believe, not of the same origin, but also, certainly, not with the same development' as Newton's method. Further improvements to the method, and its form as we know it today, were given by Thomas Simpson in his *Essays in Mathematicks* (1740). Simpson presents it as 'a new method for the solution of equations' using the 'method of fluxions', *i.e.*, derivatives. It is argued in Ypma's article that Simpson's contributions to this subject have been underestimated, and 'it would seem that the Newton–Raphson–Simpson method is a designation more nearly representing facts of history of this method which lurks inside millions of modern computer programs and is printed with Newton's name attached in so many textbooks'.

The convergence analysis of Newton's method was initiated in the first half of the twentieth century by L.V. Kantorovich. More recently, Smale, Dedieu and Shub, and others have provided significant insight into the properties of Newton's method. A full discussion of the global behaviour of the logistic equation (1.33), and other examples, will be found in P.G. Drazin, *Nonlinear Systems*, Cambridge University Press, Cambridge, 1992, particularly Chapters 1 and 3.

The secant method is also due to Newton (cf. Section 3 of Ypma's paper cited above), and is found in a collection of unpublished notes termed 'Newton's Waste Book' written around 1665.

In this chapter, we have been concerned with the iterative solution of equations for a real-valued function of a single real variable. In Chapter 4, we shall discuss the iterative solution of nonlinear systems of equations

$$\begin{aligned}
 & \text{1 : } 4(4 \quad 5, \quad 9, \quad 6 \quad , \quad \& \quad ' \\
 & \quad (\quad .@C\#.> \quad \#@ \quad .L \quad 4 \quad 4 \quad \#> \quad @ \quad 1 \\
 & \quad \quad / \quad \#> \quad 4 \\
 & \text{2} \quad 5 \quad \quad \parallel \quad \quad 1 \quad , \\
 & \quad) \quad * \quad (\quad +) \quad \quad \$ \quad , \\
 & \quad 4 \quad 4 \quad F \quad 4 \quad , \quad 4 \quad ! \quad " \quad \#.>C\#@? \quad \#@.>?4 \\
 & \text{3} \quad K \quad \$ \quad / \quad , \quad , \\
 & \quad \quad B \quad * \#D \quad \# \quad \&\#C\# \quad @. \quad 4
 \end{aligned}$$

of the form $f(z) = 0$ where f is a complex-valued function of a single complex variable z . There, corresponding to the case of $n = 2$, we shall say more about the solution of equations of the form $f(z) = 0$ where f is a complex-valued function of a single complex variable z .

This chapter has been confined to generally applicable iterative methods for the solution of a single nonlinear equation of the form $f(x) = 0$ for a real-valued function f of a single real variable. In particular, we have not discussed specialised methods for the solution of polynomial equations or the various techniques for locating the roots of polynomi-

Show that if the starting value is positive, the iteration converges to the positive solution, and if the starting value is negative it converges to the negative solution. Obtain approximate expressions for x if (i) $x = 100$ and (ii) $x = -100$, and describe the subsequent behaviour of the iteration. About how many iterations would be required to obtain the solution to six decimal digits in these two cases?

1.4 Consider the iteration

$$x_{k+1} = x_k - \frac{f(x_k)}{f'(x_k)}, \quad k = 0, 1, 2, \dots,$$

for the solution of $f(x) = 0$. Explain the connection with Newton's method, and show that (x_k) converges quadratically if x_0 is sufficiently close to the solution. Apply this method to the same example as in Example 1.7, $f(x) = e^{-x} - 2$, and verify quadratic convergence beginning from $x = 1$. Experiment with calculations beginning from $x = 10$ and from $x = -10$, and account for their behaviour.

1.5 It is sometimes said that Newton's method converges quadratically, and therefore in the successive approximations to the solution the number of correct digits doubles each time. Explain why this is not generally correct. Suppose that $f(x)$ is defined and continuous in a neighbourhood of α and that x_0 agrees with the solution α to m decimal digits; give an estimate of the number of correct decimal digits in x_1 .

Illustrate your estimate by using Newton's method to determine the positive zero of $f(x) = e^{-x} - 1.000000005$, which is close to 0.0001; use $x_0 = 0.0005$.

1.6 Suppose that $f(\alpha) = f'(\alpha) = 0$, so that f has a double root at α , and that f' is defined and continuous in a neighbourhood of α . If (x_k) is a sequence obtained by Newton's method, show that

$$x_{k+1} - \alpha = -\frac{(x_k - \alpha)^2 f''(\alpha)}{2f'(\alpha)} = -\frac{1}{2} \left(\frac{x_k - \alpha}{f'(\alpha)} \right)^2 f''(\alpha),$$

where $\alpha - m$ and $\alpha + m$ both lie between α and x_0 . Suppose, further, that $0 < m < f'(\alpha) < M$ for all x in the interval $[\alpha - m, \alpha + m]$ for some $m > 0$, where $M < 2m$; show that if x_0 lies in this interval the iteration converges to α , and that convergence is

linear, with rate $\log 2$. Verify this conclusion by finding the solution of $e^x = 1 + x$, beginning from $x = 1$.

1.7 Extend the result of the previous exercise to a case where f has a triple root at α , so that $f(\alpha) = f'(\alpha) = f''(\alpha) = 0$.

1.8 Suppose that the function f has a continuous second derivative, that $f(\alpha) = 0$, and that in the interval $[X, \beta]$, with $X < \beta$, $f'(x) > 0$ and $f''(x) < 0$. Show that the Newton iteration, starting from any x_0 in $[X, \beta]$, converges to α .

1.9 The secant method is used to determine solutions of the equation $x^3 - 1 = 0$. Starting from $x_0 = 1 + \epsilon$, $x_1 = 1 + \epsilon^2$, show that $x_2 = 1 - \epsilon + \epsilon^2$, and determine x_3 , x_4 and x_5 , neglecting terms of order ϵ^3 . Explain why, at least for sufficiently small values of ϵ , the sequence (x_k) converges to the solution 1 .

Repeat the calculation with x_0 and x_1 interchanged, so that $x_0 = 1 + \epsilon$ and $x_1 = 1 + \epsilon^2$, and show that the sequence now converges to the solution 1 .

1.10 Write the secant iteration in the form

$$x_k = \frac{x_{k-1} f(x_{k-2}) - x_{k-2} f(x_{k-1})}{f(x_{k-2}) - f(x_{k-1})}, \quad k = 1, 2, 3, \dots$$

Supposing that f has a continuous second derivative in a neighbourhood of the solution α of $f(x) = 0$, and that $f'(\alpha) > 0$ and $f''(\alpha) > 0$, define

$$(x_k, x_{k-1}) = \frac{x_k}{(x_{k-1})^2 (x_{k-2})},$$

where x_k has been expressed in terms of x_{k-1} and x_{k-2} . Find an expression for

$$(x_{k-1}) = \lim_{k \rightarrow \infty} (x_k, x_{k-1}),$$

and then determine $\lim_{k \rightarrow \infty} (x_{k-1})$. Deduce that

$$\lim_{k \rightarrow \infty} (x_k, x_{k-1}) = f'(\alpha) / 2f''(\alpha).$$

Now assume that

$$\lim_{x \rightarrow \alpha} \frac{x}{x - \alpha} = A.$$

Show that $q = 1 - 1/q = 0$, and hence that $q = -(1 + \sqrt{5})$.

Deduce finally that

$$\lim \frac{x}{x} = \frac{f(\cdot)}{2f(\cdot)}$$

Definition 2.1 The set of all $m \times n$ matrices with real entries is denoted by $\mathbb{R}^{m \times n}$. A matrix of size $n \times n$ will be called a square matrix of order n , or simply a matrix of **order n** . The **determinant** of a square matrix $A \in \mathbb{R}^{n \times n}$ is the real number $\det(A)$ defined as follows:

$$\det(A) = \sum_{\sigma \in S_n} \text{sign}(\sigma) a_{1\sigma(1)} a_{2\sigma(2)} \dots a_{n\sigma(n)}.$$

The summation is over all $n!$ permutations $(\sigma(1), \dots, \sigma(n))$ of the integers $1, 2, \dots, n$, and $\text{sign}(\sigma) = +1$ or -1 depending on whether the n -tuple $(\sigma(1), \dots, \sigma(n))$ is an even or odd permutation of $(1, 2, \dots, n)$, respectively. An even (odd) permutation is obtained by an even (odd) number of exchanges of two adjacent elements in the array $(1, 2, \dots, n)$. A matrix $A \in \mathbb{R}^{n \times n}$ is said to be **nonsingular** when its determinant $\det(A)$ is nonzero.

The **inverse matrix** A^{-1} of a nonsingular matrix $A \in \mathbb{R}^{n \times n}$ is defined as the element of $\mathbb{R}^{n \times n}$ such that $A^{-1}A = AA^{-1} = I$, where I is the $n \times n$ identity matrix

$$I = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & 1 \end{pmatrix}. \tag{2.3}$$

In order to find an explicit expression for A^{-1} in terms of the elements of the matrix A , we recall from linear algebra that, for each $i = 1, 2, \dots, n$,

$$a_{i1}A_{1j} + a_{i2}A_{2j} + \dots + a_{in}A_{nj} = \begin{cases} \det(A) & \text{if } i = j, \\ 0 & \text{if } i \neq j, \end{cases} \tag{2.4}$$

where $A_{ij} = (-1)^{i+j} \text{Cof}(a_{ij})$ and $\text{Cof}(a_{ij})$, called the **cofactor** of a_{ij} , is the determinant of the $(n-1) \times (n-1)$ matrix obtained by erasing from $A \in \mathbb{R}^{n \times n}$ row i and column j . Then, it is a trivial matter to show using (2.4) that A^{-1} has the form

$$A^{-1} = \frac{1}{\det(A)} \begin{pmatrix} A_{11} & A_{12} & \dots & A_{1n} \\ A_{21} & A_{22} & \dots & A_{2n} \\ \dots & \dots & \dots & \dots \\ A_{n1} & A_{n2} & \dots & A_{nn} \end{pmatrix}. \tag{2.5}$$

Having found an explicit formula for the matrix A^{-1} , we now multiply both sides of the equation $AX = b$ on the left by A^{-1} to deduce that

$A^{-1}(A^{-1}x) = A^{-1}x$; finally, since $A^{-1}(A^{-1}x) = (A^{-1}A)^{-1}x = I^{-1}x = x$, it follows that

$$x = A^{-1}Ax, \tag{2.6}$$

where the inverse A^{-1} of the nonsingular matrix A is given in terms of the entries of A by (2.5).

An alternative approach to the solution of the linear system $Ax = b$, called Cramer's rule, proceeds by expressing the i th entry of x as

$$x_i = D_i / D, \quad i = 1, 2, \dots, n,$$

where $D = \det(A)$, and D_i is the $n \times n$ determinant obtained by replacing the i th column of D by the entries of b . Evidently, we must require that A is nonsingular, *i.e.*, that $D = \det(A) \neq 0$. Thus, all we need to do to solve $Ax = b$ is to evaluate the $n + 1$ determinants D, D_1, \dots, D_n , each of them $n \times n$, and check that $D = \det(A)$ is nonzero; the final calculation of the elements $x_i, i = 1, 2, \dots, n$, is then trivial.

The purpose of our next example is to illustrate the application of Cramer's rule.

Example 2.1 *Suppose that we wish to solve the system of linear equations*

$$\begin{aligned} x_1 + x_2 + x_3 &= 6, \\ 2x_1 + 4x_2 + 2x_3 &= 16, \\ x_1 + 5x_2 - 4x_3 &= 3. \end{aligned}$$

The solution of such a small system can easily be found in terms of determinants, by Cramer's rule. This gives

$$x_1 = D_1 / D, \quad x_2 = D_2 / D, \quad x_3 = D_3 / D,$$

1 6 6 , B 4?D B 4#D 6 , B 4#D ,
 , 1 B 4?D A #4
 2 F \$ 0+ 9 , D4) #&* F 5 , + C K 6 #&> \$
 % 6 , C # K #& ? , B9 6 #?@.
 , #& . 6 1 4) 1
 , * * 6 , 5 N "4 1
 5 1 6 , 1 % 8
 P) E , G 1 6 1 , , 5 N B#&> D 5 6
 E G 1 , 5 1 5 4

where

$$D = \begin{pmatrix} 1 & 1 & 1 \\ 2 & 4 & 2 \\ 1 & 5 & 4 \end{pmatrix}, \quad D = \begin{pmatrix} 6 & 1 & 1 \\ 16 & 4 & 2 \\ 3 & 5 & 4 \end{pmatrix},$$

with similar expressions for D and D . To obtain the solution we therefore need to evaluate four determinants.

Now you may think that since, for A nonsingular, we have expressed the solution to $Ax = b$ in the ‘closed form’

$$x = A^{-1}b$$

and have even found a formula for A^{-1} in terms of the coefficients of A , or may simply compute the entries of x directly using Cramer’s rule, the story about the simultaneous set of linear equations (2.2) has reached its happy ending. We are sorry to disappoint you: a disturbing tale is about to unfold.

Imagine the following example: let $n = 100$, say, and suppose that you have been given all 10000 entries of a 100×100 matrix A , together with the entries of a 100-component column vector b . To avoid trivialities, let us suppose that none of the entries of A or b is equal to 0. Question: *Does the linear system $Ax = b$ have a solution? If it does, how would you find, say, the 53rd entry of the solution vector x ?* Of course, you could calculate the determinant of A and check whether it is equal to zero; if not, you could then calculate the determinant D obtained by replacing the 53rd column of A by the vector b , and the required result, by Cramer’s rule, is then the ratio of these two determinants. How much time do you think you would need to accomplish this task? An hour? A day? A month?

I imagine that you do not have a large enough sheet of paper in front of you to write down this 100×100 matrix. Let us therefore start with a somewhat simpler setting. Assume that n is any integer, $n \geq 2$, and denote by d the number of arithmetic operations that are required to calculate $\det(A)$ for A $n \times n$. For example, for a 2×2 matrix,

$$\det(A) = a_{11}a_{22} - a_{12}a_{21};$$

this evaluation requires 3 arithmetic operations – 2 multiplications and 1 subtraction – giving $d = 3$. In general, we can calculate $\det(A)$ by expanding it in the elements of its first row. This requires multiplying each of the n elements in the first row of A by a subdeterminant of size

2.1 Introduction

that we cannot by this means reduce the total by more than a factor of about n , which hardly affects our conclusion.

Our other approach to solving $Ax = b$, based on computing A^{-1} from (2.5) and writing $x = A^{-1}b$, is equally inefficient: in order to compute the inverse of an $n \times n$ matrix A using determinants, one has to calculate the determinant of A as well as n determinants of size $n - 1$ each of which then has to be divided by $\det(A)$, requiring a total of approximately

$$e n! + n e (n - 1)! + n e (n + 1)!$$

arithmetic operations, just the same as before.

The aim of this chapter is to develop alternative methods for the solution of the system of linear equations $Ax = b$. We begin by considering a classical technique, Gaussian elimination. We shall then explore its relationship to the factorisation $A = LU$ of the matrix A where L is lower triangular and U is upper triangular. It will be seen that by using the Gaussian elimination the number of arithmetic operations required to solve the linear system $Ax = b$ with an $n \times n$ matrix A is approximately $\frac{1}{2}n^3$ - a dramatic reduction from the $e(n + 1)!$ operation count associated with matrix inversion using determinants.

We conclude the chapter with a discussion of another classical idea attributed to Gauss: the least squares method for the solution of the system of linear equations $Ax = b$ where A is $n \times m$, x is the column vector of unknowns of size n and b a given column vector of size m .

2.2 Gaussian elimination

The technique for solving systems of linear algebraic equations that we shall describe in this section was developed by Carl Friedrich Gauss and was first published in his *Theoria motus corporum coelestium in sectionibus conicis solem ambientium* (1809), a major two-volume treatise on the motion of celestial bodies. Gauss was concerned with the study of

1 9 , F B* % #&&& , " / , 6 1 , " O 6
 B F 6D C * 9 6 # . >> F I O 5 F 6D
 ' #&@@ 1 9 ' 6 , 6I O 5
 E J 1 6 1 % 4 F "
 2 1 8 2# 3 ?& # 6 # #0 >? # 1604 H
 ' # 6 1 # 12 R " # 6 ' , ' G
 . ## # 140 6 || 1 B 4>D4
 3 5 : , F 4 1 , 6

the asteroid Pallas, and derived a set of six linear equations with six unknowns, also giving a systematic method for its solution.

The method proceeds by successively eliminating the elements below the diagonal of the matrix of the linear system until the matrix becomes triangular, when the solution of the system is very easy. This technique is now known under the name **Gaussian elimination**.

Before we embark on the general description of Gaussian elimination, let us illustrate its basic steps through a simple example; this is the same as Example 2.1 above, written out again for convenience.

Example 2.2 Consider the system of linear equations

$$\begin{aligned} x + x + x &= 6, \\ 2x + 4x + 2x &= 16, \\ x + 5x + 4x &= 3. \end{aligned}$$

It is convenient to rewrite this in the form $Ax = b$ where A and b are column vectors of size 3; thus,

$$\begin{pmatrix} 1 & 1 & 1 \\ 2 & 4 & 2 \\ 1 & 5 & 4 \end{pmatrix} \begin{pmatrix} x \\ x \\ x \end{pmatrix} = \begin{pmatrix} 6 \\ 16 \\ 3 \end{pmatrix}. \tag{2.9}$$

We begin by adding the first row, multiplied by -2 , to the second row, and adding the first row to the third row, giving the new system

$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & 2 & 0 \\ 0 & 6 & 3 \end{pmatrix} \begin{pmatrix} x \\ x \\ x \end{pmatrix} = \begin{pmatrix} 6 \\ 4 \\ 3 \end{pmatrix}. \tag{2.10}$$

The newly created 0 entries in the first column have been typeset in italics. Now adding the new second row, multiplied by -3 , to the third row, we find

$$\begin{pmatrix} 1 & 1 & 1 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{pmatrix} \begin{pmatrix} x \\ x \\ x \end{pmatrix} = \begin{pmatrix} 6 \\ 4 \\ 9 \end{pmatrix}, \tag{2.11}$$

1 6 " 6 " \$
 2 (4 " - B 6 4
 #@@@D , 8 1 1 6 H81 5 6
 E5 " 4 " 5 6 R 6 1 1 E5 G 6
 ' 64 % ' ' " 5 5 \$
 E ' 6 %/4 " ' O ' 6 6

which can easily be solved for the unknowns in the reverse order, beginning with $x_n = 3$.

Each of these successive row operations can be expressed as a multiplication on the left of the matrix A (in our example $n = 3$), of the system of linear equations by a transformation matrix. Writing E_{rs} for the $n \times n$ matrix whose only nonzero element is $e_{rs} = 1$, we see that the product

$$(I + \mu E_{rs})A \quad (2.12)$$

is the same as the original matrix A , except that the elements of row s , multiplied by a real number μ , have been added to the corresponding elements of row r . Here I denotes the $n \times n$ identity matrix defined by (2.3). In the elimination process we always add a multiple of an earlier row to a later row in the matrix, so that $1 \leq s < r \leq n$ in (2.12); the transformation matrix $I + \mu E_{rs}$ is therefore lower triangular in the following sense.

Definition 2.2 Let n be an integer, $n \geq 2$. The matrix $L = (l_{ij})$ is said to be **lower triangular** if $l_{ij} = 0$ for every i and j with $1 \leq i < j \leq n$. The matrix $L = (l_{ij})$ is called **unit lower triangular** if it is lower triangular, and also the diagonal elements are all equal to unity, that is $l_{ii} = 1$ for $i = 1, 2, \dots, n$.

Thus the matrix $I + \mu E_{rs}$ appearing in (2.12) is unit lower triangular if $1 \leq s < r \leq n$, and the above elimination process can be expressed by multiplying A on the left successively by the unit lower triangular matrices $I + \mu E_{rs}$ for $r = s+1, \dots, n$ and $s = 1, \dots, n-1$, with μ arbitrary; there are $n(n-1)/2$ of these matrices, one for each element of A below the diagonal (since there are n elements on the diagonal and, therefore, $1 + 2 + \dots + (n-1) = n(n-1)/2$ elements below the diagonal). The next theorem lists the technical tools which are required for proving that the resulting product is a lower triangular matrix.

Theorem 2.1 The following statements hold for any integer $n \geq 2$:

- (i) the product of two lower triangular matrices of order n is lower triangular of order n ;
- (ii) the product of two unit lower triangular matrices of order n is unit lower triangular of order n ;
- (iii) a lower triangular matrix is nonsingular if, and only if, all the

diagonal elements are nonzero; in particular, a unit lower triangular matrix is nonsingular;

- (iv) the inverse of a nonsingular lower triangular matrix of order n is lower triangular of order n ;
- (v) the inverse of a unit lower triangular matrix of order n is unit lower triangular of order n .

Proof The proofs of parts (i), (ii), (iii) and (v) are very straightforward, and are left as an exercise.

Part (iv) is proved by induction; it is easily verified for a nonsingular lower triangular matrix of order 2, using (2.5). Let $n > 2$, suppose that (iv) is true for all nonsingular lower triangular matrices of order k , with $2 \leq k < n$, and let L be a nonsingular lower triangular matrix of order $k + 1$. Both L and its inverse L^{-1} can be partitioned by their last row and column:

$$L = \begin{pmatrix} L & \mathbf{0} \\ \mathbf{x} & \alpha \end{pmatrix}, \quad L^{-1} = \begin{pmatrix} X & \mathbf{y} \\ \mathbf{z} & \beta \end{pmatrix},$$

where L is a nonsingular lower triangular matrix of order k and X is lower triangular of order k ; α, β are real numbers and $\mathbf{x}, \mathbf{y}, \mathbf{z}$ are column vectors of size k . Since the product LL^{-1} is the identity matrix of order $k + 1$, we have

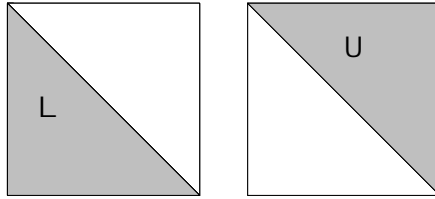
$$LX = I, \quad L\mathbf{y} = \mathbf{0}, \quad X\mathbf{z} + \beta\mathbf{z} = \mathbf{0}, \quad \alpha\beta = 1;$$

here I signifies the identity matrix of order k . Thus $X = L^{-1}$, which is lower triangular of order k by the inductive hypothesis, and $\mathbf{y} = \mathbf{0}$ given that L is nonsingular; the remaining two equations determine \mathbf{z} and β on noting that $\mathbf{z} = \mathbf{0}$ (given that L is nonsingular). This shows that L^{-1} is lower triangular of order $k + 1$, and the inductive step is complete; consequently, (iv) is true for any $n \geq 2$. □

We shall also require the concept of upper triangular matrix.

Definition 2.3 Let n be an integer, $n \geq 2$. The matrix U is said to be **upper triangular** if $u_{ij} = 0$ for every i and j with $1 \leq j < i \leq n$.

We note that results analogous to those in the preceding theorem concerning lower triangular matrices are also valid for upper triangular matrices (replacing the words ‘lower triangular’ by ‘upper triangular’ throughout).



34

5 !

The elimination process for $A \in \mathbb{R}^{n \times n}$ may now be written as follows:

$$L_1 L_2 \dots L_N A = U, \quad N = -n(n-1), \quad (2.13)$$

where $U \in \mathbb{R}^{n \times n}$ is an upper triangular matrix and each of the matrices $L_j \in \mathbb{R}^{n \times n}$, $j = 1, \dots, N$, is unit lower triangular of order n and has the form $I + \mu E_{rs}$ with $1 \leq s < r \leq n$, where I is the identity matrix of order n . That is,

$$L_1 = I + \mu_1 E_{rs}, \quad L_2 = I + \mu_2 E_{rs}, \quad \dots, \quad L_N = I + \mu_N E_{rs}.$$

It is easy to see that $E_{rs} E_{rs} = E_{rs}$, where

$$E_{rs} = \begin{cases} 1 & \text{for } r = s, \\ 0 & \text{for } r \neq s \end{cases}$$

is known as the **Kronecker delta**. Thus, for $1 \leq s < r \leq n$, the inverse of the matrix $I + \mu E_{rs}$ is the lower triangular matrix $I - \mu E_{rs}$, which corresponds to the subtraction of row s , multiplied by μ , from row r . Hence

$$A = L_1^{-1} \dots L_N^{-1} U = LU, \quad (2.14)$$

where L , as the product of a finite number of unit lower triangular matrices of order n , is itself unit lower triangular of order n by Theorem 2.1(ii); see Figure 2.1.

2.3 LU factorisation

Having seen that the Gaussian elimination process gives rise to the factorisation $A = LU$ of the matrix $A \in \mathbb{R}^{n \times n}$, $n \geq 2$, where L is unit

$$L = \begin{pmatrix} 1 & & & \\ \mu_{21} & 1 & & \\ \mu_{31} & \mu_{32} & 1 & \\ \mu_{41} & \mu_{42} & \mu_{43} & 1 \end{pmatrix}, \quad U = \begin{pmatrix} u_{11} & u_{12} & \dots & u_{1n} \\ 0 & u_{22} & \dots & u_{2n} \\ \dots & \dots & \dots & \dots \\ 0 & 0 & \dots & u_{nn} \end{pmatrix}$$

lower triangular and \mathbf{U} is upper triangular, we shall now show how to calculate the elements of \mathbf{L} and \mathbf{U} directly. Equating the elements of \mathbf{A} and \mathbf{LU}

in the formula (2.18) is zero. To investigate this possibility we use the properties of certain submatrices of A .

Definition 2.4 Suppose that A is an $n \times n$ matrix with $n \geq 2$, and let $1 \leq k < n$. The **leading principal submatrix** of order k of A is defined as the matrix A_k whose element in row i and column j is equal to the element of the matrix A in row i and column j for $1 \leq i, j \leq k$.

Armed with this definition, we can now formulate the main result of this section. It provides a sufficient condition for ensuring that the algorithm (2.18), (2.19) for calculating the entries of the matrices L and U in the LU factorisation $A = LU$ of a matrix A does not break down due to division by zero in (2.18).

Theorem 2.2 Let $n \geq 2$, and suppose that A is such that every leading principal submatrix A_k of A of order k , with $1 \leq k < n$, is nonsingular. (Note that A itself is not required to be nonsingular.) Then, A can be factorised in the form $A = LU$, where L is unit lower triangular and U is upper triangular.

Proof The proof is by induction on the order n . Let us begin by verifying the statement of the theorem for $n = 2$. We intend to show that any 2×2 matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix},$$

with $a \neq 0$, is equal to the product of a unit lower triangular matrix L of order 2 and an upper triangular matrix U of order 2; that is, we wish to establish the existence of

$$L = \begin{pmatrix} 1 & 0 \\ m & 1 \end{pmatrix}, \quad U = \begin{pmatrix} u & v \\ 0 & \end{pmatrix},$$

such that $LU = A$, where m , u , v and d are four real numbers, to be determined. Equating the product LU with A , we deduce that

$$u = a, \quad v = b, \quad mu = c, \quad mv + d = d.$$

Since $a \neq 0$ by hypothesis, the first of these equalities implies that $u \neq 0$ also; hence $m = c/u$, $v = b$, and $d = d - mv$. Thus we have shown the existence of the required matrices L and U in 2×2 and completed the proof for $n = 2$.

Now, suppose that the statement of the theorem has already been verified for matrices of order k , $2 \leq k < n$; suppose that A and all leading principal submatrices of A of order k and less are nonsingular. We mimic the proof in the case of $n = 2$ by partitioning A into blocks by the last row and column:

$$A = \begin{pmatrix} A & \mathbf{a} \\ \mathbf{b} & d \end{pmatrix}$$

where A is a nonsingular matrix (all of whose leading principal submatrices are themselves nonsingular), \mathbf{a} and \mathbf{b} are column vectors of size k , and d is a real number. According to our inductive hypothesis, there exist a unit lower triangular matrix L of order k and an upper triangular matrix U of order k such that $A = LU$. Thus we shall seek the desired unit lower triangular matrix L of order $k + 1$ and the upper triangular matrix U of order $k + 1$ in the form

$$L = \begin{pmatrix} L & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix} \quad \text{and} \quad U = \begin{pmatrix} U & \mathbf{0} \\ \mathbf{0} & d \end{pmatrix}$$

where \mathbf{a} and \mathbf{b} are column vectors of size k and d is a real number, to be determined from the requirement that the product LU be equal to the matrix A . On equating LU with A , we obtain

$$LU = A, \quad L = \begin{pmatrix} L & \mathbf{0} \\ \mathbf{0} & 1 \end{pmatrix}, \quad U = \begin{pmatrix} U & \mathbf{0} \\ \mathbf{0} & d \end{pmatrix}, \quad \mathbf{a} + d\mathbf{b} = d\mathbf{a}.$$

The first of these four equalities provides no new information. However, we can use the remaining three to determine the column vectors \mathbf{a} and \mathbf{b} and the real number d . Since L is unit lower triangular, its determinant is equal to 1; therefore L is nonsingular. This means that the second equation uniquely determines the unknown column vector \mathbf{a} . Further, since by hypothesis A is nonsingular and $A = LU$, we conclude that

$$\det(A) = \det(LU) = \det(L)\det(U) = \det(U);$$

given that $\det(A) = 0$ by the inductive hypothesis, this implies that $\det(U) = 0$ also, and therefore the third equation uniquely determines d . Having found \mathbf{a} and d , the fourth equation yields $\mathbf{b} = d^{-1}(\mathbf{a} - \mathbf{a})$. Thus we have shown the existence of the desired matrices L and U of order $k + 1$, and the inductive step is complete. \square

1) $\begin{pmatrix} 1 & & & & \\ & \ddots & & & \\ & & \ddots & & \\ & & & \ddots & \\ & & & & 1 \end{pmatrix}$! " # A $\begin{pmatrix} & & & & \\ & & & & \\ & & & & \\ & & & & \\ & & & & \end{pmatrix}$ B D A
 B D B D 5 # # 6 6 % B D A
 , 6 B#&. @C#. >&D K , G B#&. ?C#. >?D \$:

2.4 Pivoting

The aim of this section is to show that even if the matrix \mathbf{A} does not satisfy the conditions of Theorem 2.2, by permuting rows and columns it can be transformed into a new matrix $\tilde{\mathbf{A}}$ of the same size so that $\tilde{\mathbf{A}}$ admits an LU factorisation.

Example 2.3 Consider, for example, the system obtained from (2.9) by replacing the coefficient of x in the first equation by zero. Then, the leading element in the matrix \mathbf{A} is zero, the computation fails at the first step, and the LU factorisation of \mathbf{A} does not exist. However if we interchange the first two equations we obtain a new matrix $\tilde{\mathbf{A}}$ which is the same as \mathbf{A} but with the first two rows interchanged,

$$\tilde{\mathbf{A}} = \begin{pmatrix} 2 & 4 & 2 \\ 0 & 1 & 1 \\ 1 & 5 & 4 \end{pmatrix}. \quad (2.20)$$

Since the leading principal submatrices of order 1 and 2 of $\tilde{\mathbf{A}}$ are non-singular, by Theorem 2.2 the matrix $\tilde{\mathbf{A}}$ now has the required LU factorisation, which is easily computed.

A computation which fails when an element is exactly zero is also likely to run into difficulties when that element is nonzero but of very small absolute value; the problem stems from the presence of rounding errors. The basic operation in the elimination process consists of multiplying the elements of one row of the matrix by a scalar μ , and adding to the elements of another row. The multiplication operation will always introduce a rounding error, so the elements which are multiplied by μ will already contain a rounding error from operations with earlier rows of the matrix; these errors will therefore themselves be multiplied by μ before adding to the new row. The errors will be magnified if $|\mu| > 1$, and will be greatly magnified if $|\mu| \gg 1$.

The accumulation of rounding errors alluded to in the previous paragraph can be alleviated by permuting the rows of the matrix. Thus, at each stage of the elimination process we interchange two rows, if necessary, so that the largest element in the current column lies on the diagonal. This process is known as **pivoting**. Clearly, when pivoting is performed none of the multipliers μ have absolute value greater than unity. The process is easily formalised by introducing permutation matrices. This leads us to our next definition.

Definition 2.5 Suppose that $n \geq 2$. A matrix $P \in \mathbb{R}^{n \times n}$ in which every element is either 0 or 1, and whose every row and every column contain exactly one nonzero element, is called a **permutation matrix**.

Example 2.4 Here are three of the possible $3!$ permutation matrices in $\mathbb{R}^{3 \times 3}$:

$$\begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \begin{pmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \begin{pmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix}.$$

The proof of our next result is elementary and is left to the reader.

Lemma 2.1 Let $n \geq 2$ and suppose that $P \in \mathbb{R}^{n \times n}$ is a permutation matrix. Then, the following statements hold:

- (i) given that I is the identity matrix of order n , the matrix P can be obtained from I by permuting rows;
- (ii) if $Q \in \mathbb{R}^{n \times n}$ is another permutation matrix, then the products PQ and QP are also permutation matrices;
- (iii) let $P_{rs} \in \mathbb{R}^{n \times n}$ denote the **interchange matrix**, obtained from the identity matrix $I \in \mathbb{R}^{n \times n}$ by interchanging rows r and s ; any interchange matrix is a permutation matrix; moreover, any permutation matrix of order n can be written as a product of interchange matrices of order n ;
- (iv) the determinant of a permutation matrix $P \in \mathbb{R}^{n \times n}$ is equal to 1 or -1 , depending on whether P is obtained from the identity matrix of order n by an even or odd number of permutations of rows, respectively; in particular, a permutation matrix is nonsingular.

Now we are ready to prove the next theorem.

Theorem 2.3 Let $n \geq 2$ and $A \in \mathbb{R}^{n \times n}$. There exist a permutation matrix P , a unit lower triangular matrix L , and an upper triangular matrix U , all three in $\mathbb{R}^{n \times n}$, such that

$$PA = LU. \tag{2.21}$$

Proof The proof is by induction on the order n . Let $n = 2$ and consider the matrix

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}.$$

If $a = 0$, the proof follows from Theorem 2.2 with P taken as the 2×2 identity matrix. If $a = 0$ but $c = 0$, we take

$$P = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

and write

$$PA = \begin{pmatrix} c & d \\ 0 & b \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} c & d \\ 0 & b \end{pmatrix} = LU.$$

If $a = 0$ and $c = 0$, the result trivially follows by writing

$$\begin{pmatrix} 0 & b \\ 0 & d \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 0 & b \\ 0 & d \end{pmatrix} = LU$$

and taking P as the 2×2 identity matrix. That completes the proof for $n = 2$.

Now, suppose that A is $n \times n$ and assume that the theorem holds for every matrix of order k with $2 \leq k < n$. We begin by locating the element in the first column of A which has the largest absolute value, or any one of them if there is more than one such element, and interchange rows if required; if the largest element is in row r we interchange rows 1 and r . We then partition the new matrix according to the first row and column, writing

$$P A = \begin{pmatrix} l & \mathbf{0} \\ \mathbf{B} & \mathbf{C} \end{pmatrix} \tag{2.22}$$

where l is the element of largest absolute value in the first column, \mathbf{B}, \mathbf{C} are $(n-1) \times (n-1)$ matrices, and \mathbf{b}, \mathbf{c} and \mathbf{C} are column vectors of size $n-1$, with \mathbf{b}, \mathbf{c} and \mathbf{C} to be determined. Writing out the product we find that

$$\begin{aligned} \mathbf{b} &= \mathbf{0}, \\ \mathbf{c} &= \mathbf{0}, \\ \mathbf{C} &= \mathbf{B}^{-1} \mathbf{c}. \end{aligned} \tag{2.23}$$

If $l = 0$, then the first column of A consists entirely of zeros ($\mathbf{b} = \mathbf{0}$); in this case we can evidently choose $\mathbf{b} = \mathbf{0}$, $\mathbf{c} = \mathbf{0}$ and $\mathbf{C} = \mathbf{B}$. Suppose now that $l \neq 0$; then $\mathbf{b} = (1/l) \mathbf{c}$, so that all the elements of \mathbf{b} have absolute value less than or equal to unity, since l is the largest in absolute value element in the first column. By the inductive hypothesis we can now write

$$P C = L U, \tag{2.24}$$

where P , L , U are $n \times n$ matrices, P is a permutation matrix, L is unit lower triangular, and U is upper triangular. Hence, by (2.23),

$$PA = \begin{pmatrix} 1 & \mathbf{0} & & \\ \mathbf{0} & P & & \\ & & 1 & \mathbf{0} \\ & & & L \\ & & & & \mathbf{0} & U \end{pmatrix} \quad (2.25)$$

since $PP^{-1} = I$. Now, defining the permutation matrix P by

$$P = \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & P \end{pmatrix}, \quad (2.26)$$

we obtain

$$PA = \begin{pmatrix} 1 & \mathbf{0} & & \\ P & L & & \\ & & \mathbf{0} & U \end{pmatrix}, \quad (2.27)$$

which is the required factorisation of A . This completes the inductive step. The theorem therefore holds for every matrix of order $n \geq 2$. \square

The proof of this theorem also contains an algorithm for constructing the permutation matrix P , and the matrices L and U . The permutation matrix is conveniently described by specifying the sequence of interchanges: given the $n - 1$ integers p_1, p_2, \dots, p_{n-1} , the matrix P is the product of the permutation matrices which interchange rows 1 and p_1 , 2 and p_2 , and so on.

2.5 Solution of systems of equations

Consider the linear system $Ax = b$ where A is $n \times n$ and b and x are column vectors of size n . According to Theorem 2.3 there exist a permutation matrix P , a unit lower triangular matrix L and an upper triangular matrix U such that $PA = LU$. Having obtained the LU factorisation of the matrix PA , the solution of the system of linear equations $Ax = b$ is straightforward: multiplying both sides of $Ax = b$ on the left by the permutation matrix P , we obtain that

$$PAx = Pb; \quad (2.28)$$

equivalently, $LUx = Pb$. On defining $y = Ux$ we can rewrite (2.28) as the following coupled set of linear equations:

$$Ly = Pb, \quad Uy = x. \quad (2.29)$$

Assuming that the matrix P and the LU factorisation of PA are already known, there are three stages to the calculation of x :

() First we apply the sequence of permutations to the vector y , to produce $P^{-1}y$;

()* We then solve the lower triangular system $Ly = P^{-1}y$, calculating the elements in the order y_1, y_2, \dots, y_n ;

()- Finally the required solution x is obtained from the upper triangular system $Ux = P^{-1}y$, calculating the elements of x in the reverse order, x_n, x_{n-1}, \dots, x_1 .

(,) will break down if any of the diagonal elements of U are zero, but if this happens the matrix A is singular.

The next section is devoted to assessing the amount of computational work for this algorithm.

2.6 Computational work

In this section we shall show that the work involved in factorising an $n \times n$ matrix in the form $A = LU$ is proportional to n^3 . An estimate of the amount of computational work of this kind is important in deciding in advance how long a calculation would take for a very large matrix, and is also useful in comparing different methods for the solution of a given problem. For example, in the next chapter we shall derive a method for solving a system of equations with a symmetric positive definite matrix; that method requires only half the amount of work involved in the standard LU factorisation algorithm which takes no account of symmetry.

Accurate estimates of the time taken by a computation are very complicated and require some detailed knowledge of the computer being used. The estimates which we shall give are simple but crude; they are normally good enough for the types of comparisons we have just mentioned.

We see from (2.18) that the calculation of l_{ij} requires $j - 1$ multiplications, $j - 2$ additions, 1 subtraction and 1 division, a total of $2j - 1$ operations. In the same way, (2.19) shows that the calculation of u_{ik} requires $2i - 2$ operations. Recalling that, for any integer $k \geq 2$,

$$1 + 2 + \dots + k = \frac{1}{2}k(k + 1) \quad \text{and} \quad 1 + 2 + \dots + k = \frac{1}{2}k(k + 1)(2k + 1),$$

we then deduce that the total number of operations involved in the LU

factorisation is

$$(2j-1) + \dots + 2(i-1) = -n(n-1)(4n+1).$$

It is enough to say that the number of multiplications required is about $-n^3 - n$, for moderately large values of n .

Having constructed the factorisation we can now count the number of operations required to compute the vectors y_i and u_i in (2.29). Given the vector P_i , the elements of y_i are obtained from

$$y_i = (P_i)^{-1}, \quad y_i = (P_i)^{-1} \mid y_i, \quad i = 2, 3, \dots, n, \quad (2.30)$$

which requires $2i-2$ operations. Summing over i this gives a total of $n(n-1)$. The calculation of the elements of u_i is similar:

$$x = \frac{1}{u} y \quad u = x, \quad i = 1, 2, \dots, n. \quad (2.31)$$

This requires $2(n-i)+1$ operations, giving a total of n .

The total number of operations involved in the solution of the system of equations is therefore approximately $-n^3 - n$ for the factorisation, followed by $n(n-1)+n = 2n^2 - n$ for the solution of the two triangular systems, that is, approximately $-n^3 + -n$, ignoring terms of size (n) .

We often need to solve a number of systems of this kind, all with different right-hand sides, but with the same matrix. We then need only factorise the matrix once, and the total number of multiplications required for k right-hand sides becomes approximately $-n^3 + 2kn^2 - n$. When k is fairly large it might appear that it would be more efficient to form the inverse matrix A^{-1} , and then multiply each right-hand side by the inverse; but we shall show that it is not so.

To form the inverse matrix we first factorise the matrix A , and then solve n systems, with the right-hand sides being the vectors which constitute the columns of the identity matrix. Because these right-hand sides have a special form, there is the possibility of saving some work; some careful counting shows that the total can be reduced from $-n^3 + 2n^2 = -n^3$ to an approximate total of $2n^2$ operations. It is easy to see that the operation of multiplying a vector by the inverse matrix requires $n(2n-1)$ operations; hence the whole computation of first constructing the inverse matrix, and then multiplying each right-hand side by the inverse, requires a total of $2n^2 + 2kn^2$ multiplications (ignoring terms of size

(n)). This is always greater than the previous value $-n + 2k - n$, whether k is small or large. The most efficient way of solving this problem is to construct and save the L and U factors of A , rather than to form the inverse of A .

2.7 Norms and condition numbers

The analysis of the effects of rounding error on solutions of systems of linear equations requires an appropriate measure. This is provided by the concept of **norm** defined below. In order to motivate the axioms of norm stated in Definition 2.6, we note that the set of real numbers is a linear space, and that the **absolute value** function

$$|v| = \begin{cases} v & \text{if } v \geq 0, \\ -v & \text{if } v < 0 \end{cases}$$

has the following properties:

$$\begin{aligned} |v| &\geq 0 \text{ for any } v, \text{ and } |v| = 0 \text{ if, and only if, } v = 0; \\ |v| &= |-v| \text{ for all } v \text{ and all } v; \\ |u+v| &\leq |u| + |v| \text{ for all } u \text{ and } v \text{ in } \mathbb{R}. \end{aligned}$$

The absolute value $|v|$ of a real number v measures the distance between v and 0 (the zero element of the linear space \mathbb{R}). Our next definition aims to generalise this idea to an arbitrary linear space over the field of real numbers: even though the discussion in the present chapter is confined to finite-dimensional linear spaces of vectors ($V = \mathbb{R}^n$) and square matrices ($M = \mathbb{R}^{n \times n}$), norms over other linear spaces, including infinite-dimensional function spaces, will appear elsewhere in the text (see Chapters 8, 9, 11 and 14).

Definition 2.6 Suppose that V is a linear space over the field of real numbers. The nonnegative real-valued function $\| \cdot \|$ is said to be a **norm** on the space V provided that it satisfies the following axioms:

$$\begin{aligned} \|v\| &\geq 0 \text{ if, and only if, } v = 0 \text{ in } V; \\ \|v\| &= \|-v\| \text{ for all } v \text{ and all } v \text{ in } V; \\ \|u+v\| &\leq \|u\| + \|v\| \text{ for all } u \text{ and } v \text{ in } V \text{ (the triangle inequality)}. \end{aligned}$$

A linear space V , equipped with a norm, is called a **normed linear space**.

Remark 2.1 If V is a linear space over the field of complex numbers, then $\| \cdot \|$ in the second axiom of Definition 2.6 should be replaced by $\| \cdot \|$, with $\| \cdot \|$ signifying the modulus of \cdot .

Any norm on the linear space \mathbb{R}^n will be called a **vector norm**. Three vector norms are in common use in numerical linear algebra: the 1-norm $\| \cdot \|_1$, the 2-norm (or Euclidean norm) $\| \cdot \|_2$, and the ∞ -norm $\| \cdot \|_\infty$; these are defined below.

Definition 2.7 The 1-norm of the vector $v = (v_1, \dots, v_n)$ is defined by

$$\|v\|_1 = |v_1| + \dots + |v_n|. \tag{2.32}$$

Definition 2.8 The 2-norm of the vector $v = (v_1, \dots, v_n)$ is defined by $\|v\|_2 = \sqrt{v_1^2 + \dots + v_n^2}$. In other words,

$$\|v\|_2 = \sqrt{\sum_{i=1}^n v_i^2}. \tag{2.33}$$

Definition 2.9 The ∞ -norm of the vector $v = (v_1, \dots, v_n)$ is defined by

$$\|v\|_\infty = \max_{1 \leq i \leq n} |v_i|. \tag{2.34}$$

When $n = 1$, each of these norms collapses to the absolute value, $\|x\|_1 = \|x\|_2 = \|x\|_\infty = |x|$, the simplest example of a norm on \mathbb{R} .

It is easy to show that $\| \cdot \|_1$ and $\| \cdot \|_\infty$ obey all axioms of a norm. For the 2-norm the first two axioms are still trivial to verify; to show that the triangle inequality is satisfied by the 2-norm requires use of the Cauchy–Schwarz inequality.

Lemma 2.2 (Cauchy–Schwarz inequality)

$$\|u + v\|_2^2 = \|u\|_2^2 + \|v\|_2^2 + 2(u, v). \tag{2.35}$$

1 % \$: , 6 B # % #&. @ 9 , C * ' 6 #.>& , 8
 B D 9 , D 5 6 E, , 6 6
 1, 64 O 81 , , 1 1 1 , 1 , 6 64 6
 2 , O % , + B > K 6 #. * O 1 F 6
 B D C * 5 # @ # F 6D , ,
 , , 1 1' , , #. @ 4 H ,
 , , 1 6 , (, 6 9 , 1 4 \$

Theorem 2.4 (Young’s inequality) Let $p, q > 1$, $(1/p) + (1/q) = 1$. Then, for any two nonnegative real numbers a and b ,

$$ab \leq \frac{a^p}{p} + \frac{b^q}{q}.$$

Proof If either $a = 0$ or $b = 0$ the inequality holds trivially. Let us therefore suppose that $a > 0$ and $b > 0$. We recall that a function $f(x)$ is said to be **convex** if

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

for all $\lambda \in [0, 1]$, and all x and y in \mathbb{R} ; i.e., for any x and y in \mathbb{R} the graph of the function f between the points $(x, f(x))$ and $(y, f(y))$ lies below the chord that connects these two points. Note that the function x^e is convex. Therefore, with $\lambda = 1/p$ and $1 - \lambda = 1/q$, we get that

$$ab = e^{\lambda \ln a + (1-\lambda) \ln b} = e^{\frac{1}{p} \ln a + \frac{1}{q} \ln b} = \frac{a}{p} + \frac{b}{q},$$

and the proof is complete. (When $p = q = 2$ the proof is trivial: as $(a - b)^2 \geq 0$ also $2ab \leq a^2 + b^2$, and hence the required result.) \square

The next step is to establish Hölder’s inequality; it is a generalisation of the Cauchy–Schwarz inequality.

Theorem 2.5 (Hölder’s inequality) Let $p, q > 1$, $(1/p) + (1/q) = 1$. Then, for any $u_i, v_i \geq 0$ and $i = 1, \dots, n$, we have

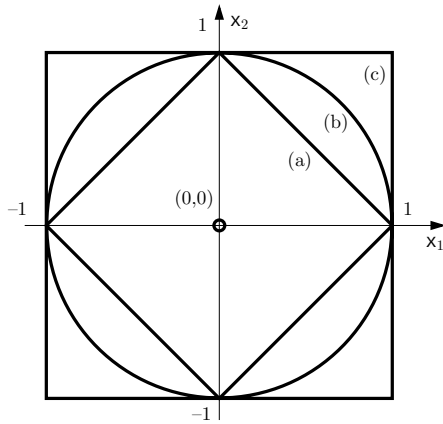
$$\sum_{i=1}^n u_i v_i \leq \left(\sum_{i=1}^n u_i^p \right)^{1/p} \left(\sum_{i=1}^n v_i^q \right)^{1/q}.$$

Proof If either $\sum u_i^p = 0$ or $\sum v_i^q = 0$ the inequality holds trivially. Let us therefore suppose that $\sum u_i^p > 0$ and $\sum v_i^q > 0$, and consider the vectors \tilde{u} and \tilde{v} in \mathbb{R}^n with components $\tilde{u}_i = u_i / \left(\sum u_i^p \right)^{1/p}$ and $\tilde{v}_i = v_i / \left(\sum v_i^q \right)^{1/q}$, respectively, $i = 1, 2, \dots, n$. By Young’s inequality,

$$\tilde{u}_i \tilde{v}_i \leq \tilde{u}_i^p \tilde{v}_i^q + \frac{1}{p} \tilde{u}_i + \frac{1}{q} \tilde{v}_i = \frac{1}{p} + \frac{1}{q} = 1.$$

Inserting the defining expressions for \tilde{u} and \tilde{v} into the left-most expression in this chain, the result follows. \square

¹ H : OI B / , #.>@ F 6 C @ % #*&
 : + F 6D , 6L , , 11 ,
 4 OI , 5 , G 6
 1 #. " , 5 , 19 4



64 " 0 # 2
 * + " * , + " * + "

the unit sphere in a normed linear space V , with norm $\|\cdot\|$, is defined as the set $S = \{x \in V : \|x\| = 1\}$. It can be seen from Figure 2.2 that

$$S = \{x \in V : \|x\| = 1\} = \{x \in V : \|x\| = 1\}.$$

We leave it to the reader as an exercise to show that analogous inclusions hold in V for any $n \geq 1$. (See Exercise 8.)

The unit sphere in a normed linear space V with norm $\|\cdot\|$ is the boundary of the closed unit ball $\bar{B}(0)$ centred at 0 defined by

$$\bar{B}(0) = \{x \in V : \|x\| \leq 1\}.$$

Analogously, the open unit ball centred at 0 is defined by

$$B(0) = \{x \in V : \|x\| < 1\}.$$

More generally, for $r > 0$ and $x_0 \in V$,

$$\bar{B}(x_0, r) = \{x \in V : \|x - x_0\| \leq r\}$$

is the **closed ball** of radius r centred at x_0 ; analogously,

$$B(x_0, r) = \{x \in V : \|x - x_0\| < r\}$$

is the **open ball** of radius r centred at x_0 .

Any norm on the linear space M_n of $n \times n$ matrices with real entries will be referred to as a **matrix norm**. In particular, we shall now

consider matrix norms which are induced by vector norms in a sense that will be made precise in the next definition.

Definition 2.10 Given any norm $\|\cdot\|$ on the space $M_n(\mathbb{R})$ of n -dimensional vectors with real entries, the **subordinate matrix norm** on the space $M_n(\mathbb{R})$ of $n \times n$ matrices with real entries is defined by

$$\|A\| = \max_{\|x\|=1} \|Ax\|. \quad (2.38)$$

In (2.38) we used $\|x\|$ to denote $\|x\|$, where, for sets A and B , $A \subseteq B = \{x \in A : x \in B\}$.

Remark 2.4 Let $M_n(\mathbb{C})$ denote the linear space of $n \times n$ matrices with complex entries over the field \mathbb{C} of complex numbers. Given any norm $\|\cdot\|$ on the linear space $M_n(\mathbb{C})$, the **subordinate matrix norm** on $M_n(\mathbb{C})$ is defined by

$$\|A\| = \max_{\|x\|=1} \|Ax\|,$$

where $\|x\| = \|x\|$.

It is easy to show that a subordinate matrix norm satisfies the axioms of norm listed in Definition 2.6; the details are left as an exercise. Definition 2.10 implies that, for $A \in M_n(\mathbb{R})$,

$$\|A\| = \|A\|, \quad \text{for all } A \in M_n(\mathbb{R}).$$

In a relation like this any vector norm may be used, but of course it is necessary to use the same norm throughout. It follows from Definition 2.10 that, in any subordinate matrix norm on $M_n(\mathbb{R})$,

$$\|I\| = 1$$

where I is the $n \times n$ identity matrix.

Given any vector norm $\|\cdot\|$ in \mathbb{R}^n , it is a trivial matter to evaluate each of the three norms $\|x\|$, $\|y\|$, $\|z\|$; however, it is not yet obvious how one can calculate the corresponding subordinate matrix norm of a given matrix A in $M_n(\mathbb{R})$. Definition 2.10 is unhelpful in this respect: calculating $\|A\|$ via (2.38) would involve the unpleasant task of maximising the function $\|Ax\|$ over $\|x\|=1$ (or, equivalently, maximising $\|Ax\|$ over the unit sphere $S^{n-1} = \{x \in \mathbb{R}^n : \|x\|=1\}$). This difficulty is resolved by the following three theorems.

Theorem 2.7 *The matrix norm subordinate to the vector norm can be expressed, for an $n \times n$ matrix $A = (a_{ij})$, as*

$$\|A\|_\infty = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}|. \tag{2.39}$$

This result is often loosely expressed by saying that the ∞ -norm of a matrix is its largest row-sum.

Proof Given an arbitrary vector v in \mathbb{R}^n , write $K = \|v\|_\infty$, so that $|v_j| \leq K$ for $j = 1, 2, \dots, n$. Then,

$$(Av)_i = \sum_{j=1}^n a_{ij} v_j \leq \sum_{j=1}^n |a_{ij}| |v_j| \leq K \sum_{j=1}^n |a_{ij}|, \quad i = 1, 2, \dots, n.$$

Now we define

$$C = \max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| \tag{2.40}$$

and note that

$$\frac{\|Av\|_\infty}{\|v\|_\infty} = \frac{\max_{1 \leq i \leq n} |(Av)_i|}{K} = \frac{\max_{1 \leq i \leq n} \sum_{j=1}^n |a_{ij}| |v_j|}{K} \leq C.$$

Hence, $\|A\|_\infty \leq C$.

Next we show that $\|A\|_\infty \geq C$. To do so, we take v to be a vector each of whose entries is ± 1 , with the choice of sign to be made clear below. In the definition of C , equation (2.40), let m be the value of i for which the maximum is attained, or any one of the values if there is more than one. Then, in the vector v we give the element v_m the same sign as that of a_{mj} ; if a_{mj} happens to be zero, the choice of the sign of v_m is irrelevant. With this definition of v we see at once that

$$\|Av\|_\infty = \sum_{j=1}^n |a_{mj}| |v_j| = \sum_{j=1}^n |a_{mj}| = C.$$

As $\|v\|_\infty = 1$, it follows that

$$\|A\|_\infty \geq C,$$

which means that $\|A\|_\infty = C$. Hence $\|A\|_\infty = C$, as required. \square

Theorem 2.8 *The matrix norm subordinate to the vector norm can be expressed, for an $n \times n$ matrix $A = (a_{ij})$, as*

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

This is often loosely expressed by saying that the 1-norm of a matrix is its largest column-sum. The proof of this theorem is very similar to that of the previous one, and is left as an exercise (see Exercise 7). Note that Theorems 2.7 and 2.8 mean that the 1-norm of a matrix $A = (a_{ij})$ is the ∞ -norm of the transpose $A^T = (a_{ji})$ of the matrix.

Before we state a characterisation of the subordinate matrix 2-norm, we recall the following definition from linear algebra.

Definition 2.11 *Suppose that $\lambda \in \mathbb{C}$. A complex number λ , for which the set of linear equations*

$$(A - \lambda I)x = 0$$

*has a nontrivial solution $x \neq 0$, is called an **eigenvalue** of A ; the associated solution x is called an **eigenvector** of A (corresponding to λ).*

Now we are ready to state our result.

Theorem 2.9 *Let $A \in \mathbb{R}^{n \times n}$ and denote the eigenvalues of the matrix $B = A + A^T$ by $\lambda_i, i = 1, 2, \dots, n$. Then,*

$$\|A\|_2 = \max_{1 \leq i \leq n} |\lambda_i|.$$

Proof Note first that the matrix B is symmetric, i.e., $B = B^T$; therefore all of its eigenvalues are real and the associated eigenvectors belong to \mathbb{R}^n . (You may wish to prove this: consult the proof of Theorem 3.1, part (ii), for a hint.) Moreover, all eigenvalues of B are nonnegative, since if x is an eigenvector of B and λ is the associated eigenvalue, then

$$A^T A x = B x = \lambda x$$

and therefore

$$\lambda \|x\|_2^2 = x^T A^T A x = (Ax)^T (Ax) = \|Ax\|_2^2 \geq 0.$$

Suppose that the vectors $x_i, i = 1, 2, \dots, n$, are eigenvectors of B corresponding to the eigenvalues $\lambda_i, i = 1, 2, \dots, n$. Since B is symmetric

we may assume that the vectors \mathbf{e}_i are orthogonal, *i.e.*, $\mathbf{e}_i^T \mathbf{e}_j = 0$ for $i \neq j$, and we can normalise them so that $\mathbf{e}_i^T \mathbf{e}_i = 1$ for $i = 1, 2, \dots, n$. Now choose an arbitrary vector \mathbf{c} in \mathbb{R}^n and express it as a linear combination of the vectors \mathbf{e}_i , $i = 1, 2, \dots, n$:

$$\mathbf{c} = c_1 \mathbf{e}_1 + c_2 \mathbf{e}_2 + \dots + c_n \mathbf{e}_n.$$

Then,

$$\mathbf{B} \mathbf{c} = c_1 \mathbf{B} \mathbf{e}_1 + c_2 \mathbf{B} \mathbf{e}_2 + \dots + c_n \mathbf{B} \mathbf{e}_n.$$

We may assume, without loss of generality, that

$$(0 \leq c_1 \leq c_2 \leq \dots \leq c_n).$$

Using the orthonormality of the vectors \mathbf{e}_i , $i = 1, 2, \dots, n$, we get that

$$\begin{aligned} \mathbf{A} \mathbf{B} \mathbf{c} &= \mathbf{A} (\mathbf{B} \mathbf{c}) = \mathbf{B} \mathbf{c} \\ &= c_1 \mathbf{B} \mathbf{e}_1 + c_2 \mathbf{B} \mathbf{e}_2 + \dots + c_n \mathbf{B} \mathbf{e}_n \\ &= (c_1 \mathbf{B} \mathbf{e}_1 + c_2 \mathbf{B} \mathbf{e}_2 + \dots + c_n \mathbf{B} \mathbf{e}_n) \end{aligned} \tag{2.41}$$

for any vector \mathbf{c} . Hence $\mathbf{A} \mathbf{B} \mathbf{c} = \mathbf{B} \mathbf{c}$. To prove equality we simply choose $\mathbf{c} = \mathbf{e}_i$ in (2.41), so that $c_i = 1$, $c_j = 0$ for $j \neq i$ and $c = 1$. \square

The square roots of the (nonnegative) eigenvalues of $\mathbf{A}^T \mathbf{A}$ are referred to as the **singular values** of \mathbf{A} . Thus we have shown that the 2-norm of a matrix \mathbf{A} is equal to the largest singular value of \mathbf{A} .

If the matrix \mathbf{A} is symmetric, then $\mathbf{B} = \mathbf{A}^T \mathbf{A} = \mathbf{A} \mathbf{A}^T$, and the eigenvalues of \mathbf{B} are just the squares of the eigenvalues of \mathbf{A} . In this special case the 2-norm of \mathbf{A} is the largest of the absolute values of its eigenvalues.

Theorem 2.10 *Given that $\|\cdot\|$ is a subordinate matrix norm on $\mathbb{R}^n \times \mathbb{R}^n$,*

$$\|\mathbf{A} \mathbf{B}\| \leq \|\mathbf{A}\| \|\mathbf{B}\|$$

for any two matrices \mathbf{A} and \mathbf{B} in $\mathbb{R}^n \times \mathbb{R}^n$.

Proof From the definition of subordinate matrix norm,

$$\|\mathbf{A} \mathbf{B}\| = \max_{\|\mathbf{x}\|=1} \|\mathbf{A} \mathbf{B} \mathbf{x}\|.$$

As

$$\|\mathbf{A} \mathbf{B} \mathbf{x}\| \leq \|\mathbf{A}\| \|\mathbf{B} \mathbf{x}\|$$



while if $D = [0, 1]$, then $\text{Cond}(f) = +\infty$. Indeed, in the latter case, perturbing $x = 0$ to $x = \delta$, $0 < \delta < 1$, leads to a perturbation of the function value $f(0) = 0$ to $f(\delta) = \delta = \delta^{-1} \cdot \delta$: a magnification by a factor δ^{-1} in comparison with the size of the perturbation in x .

When $f(y) - f(x) \approx f'(x)(y - x)$ exhibits large variation as (x, y) ranges through D , it is more helpful to consider a finer, local measure of conditioning, the **absolute local condition number**, at $x \in D$, of the function f , defined by

$$\text{Cond}(f) = \sup_{\|x\| \leq \delta} \frac{f(x + \delta) - f(x)}{\delta} \cdot \frac{1}{\|f'(x)\|}. \tag{2.43}$$

Example 2.6 Let us consider the function $f: x \in D \rightarrow \bar{x}$, defined on the interval $D = (0, \infty)$. The absolute local condition number of f at $x \in D$ is $\text{Cond}(f) = 1/(2\bar{x})$. Clearly, $\lim_{x \rightarrow \infty} \text{Cond}(f) = 0$, $\lim_{x \rightarrow 0^+} \text{Cond}(f) = +\infty$.

Although the definitions (2.42) and (2.43) seem intuitive, they are not always satisfactory from the practical point of view since they depend on the magnitudes of $f(x)$ and x . A more convenient definition of conditioning is arrived at by rescaling (2.43) by the norms of $f(x)$ and x . This leads us to the notion of **relative local condition number**

$$\text{cond}(f) = \sup_{\|x\| \leq \delta} \frac{f(x + \delta) - f(x)}{\|x\|} \cdot \frac{\|x\|}{\|f(x)\|},$$

where it is implicitly assumed that $x \neq 0$ and $f(x) \neq 0$. The next example highlights the difference between the absolute local condition number and the relative local condition number of f .

Example 2.7 Let us consider the function $f: x \in D \rightarrow \bar{x}$, defined on the interval $D = (0, \infty)$. Recall from the preceding example that the absolute local condition number of f at $x \in D$ approaches $+\infty$ as x tends to zero. In contrast with this, the relative local condition number of f is $\text{cond}(f) = 1/2$ for all $x \in D$.

You may also wish to ponder the following, seemingly paradoxical, observation: $\lim_{x \rightarrow 0^+} \text{cond}(\sin) = 1$ and $\lim_{x \rightarrow 0^+} \text{cond}_-(\sin) = +\infty$, even though $\sin 0 = \sin \pi = 0$ and $\text{Cond}(\sin) = \text{Cond}_-(\sin) = 1$.

Since the present section is concerned with the solution of the linear system $Ax = b$, where $A \in \mathbb{R}^{n \times n}$ is nonsingular and $b \in \mathbb{R}^n$, let us

consider the relative local condition number of the mapping

$$A^{-1} : \quad A^{-1}$$

at $x = x_0$. We suppose that X has been equipped with a vector norm $\| \cdot \|$ and, since there is no danger of confusion, we denote the associated subordinate matrix norm by $\| \cdot \|$

There is a condition number for each norm; for example, if we use the 2-norm, then $\kappa_2(A) = \|A\|_2 \|A^{-1}\|_2$, and so on. Indeed, the size of the condition number of a matrix $A \in \mathbb{R}^{n \times n}$ is strongly dependent on the choice of the norm in (2.46). In order to illustrate the last point, let us consider the unit lower triangular matrix $A \in \mathbb{R}^{n \times n}$ defined by

$$A = \begin{pmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & 0 & \dots & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 1 \end{pmatrix}, \tag{2.46}$$

and note that its inverse is

$$A^{-1} = \begin{pmatrix} 1 & 0 & 0 & 0 & \dots & 0 \\ 1 & 1 & 0 & 0 & \dots & 0 \\ 1 & 0 & 1 & 0 & \dots & 0 \\ 1 & 0 & 0 & 1 & \dots & 0 \\ \dots & \dots & \dots & \dots & \dots & \dots \\ 1 & 0 & 0 & 0 & \dots & 1 \end{pmatrix}.$$

Since

$$\|A\|_2 = n \quad \text{and} \quad \|A^{-1}\|_2 = n,$$

it follows that $\kappa_2(A) = n^2$. On the other hand,

$$\|A\|_1 = 2 \quad \text{and} \quad \|A^{-1}\|_1 = 2.$$

so that $\kappa_1(A) = 4 \|n\|_1 = \kappa_1(A)$ when $n \geq 1$. (A question for the curious: how does the condition number $\kappa_2(A)$ of the matrix A in (2.46) depend on the size n of A ? See Exercise 11.)

It is left as an exercise to show that for a symmetric matrix A (i.e., when $A^T = A$), the 2-norm condition number $\kappa_2(A)$ is the ratio of the largest of the absolute values of the eigenvalues of A to the smallest of the absolute values of the eigenvalues (see Exercise 9).

$$A = \begin{matrix} 8 \\ =1 \\ =1 \end{matrix}$$

$$A^{-1} = \begin{matrix} 8 \\ =1 \\ =1 \end{matrix}$$

We can now assess the sensitivity of the solution of the system $Ax = b$ to changes in the right-hand side vector b .

Theorem 2.11 Suppose that A is a nonsingular matrix, b is a vector, δb is a vector, and $A(\delta x) = \delta b$, with $\|\delta b\| \leq \epsilon$. Then, $\|\delta x\| \leq \kappa(A)\epsilon$ and

$$\frac{\|\delta x\|}{\|x\|} \leq \kappa(A) \frac{\|\delta b\|}{\|b\|}.$$

Proof Evidently,

$$\delta x = A^{-1} \delta b \quad \text{and} \quad \| \delta x \| = \| A^{-1} \delta b \| \leq \| A^{-1} \| \|\delta b\|.$$

As $\|b\| > 0$ and A is nonsingular, the first of these implies that $\|x\| > 0$. Further,

$$\|A\| \|x\| \leq \|b\| \quad \text{and} \quad \|A^{-1}\| \|b\| \leq \|x\|.$$

The result follows immediately by multiplying these inequalities. \square

Owing to the effect of rounding errors during the calculation, the numerical solution of $Ax = b$ will not be exact. The numerical solution may be written $x + \delta x$, and we shall usually find that this vector satisfies the equation $A(x + \delta x) = b + \delta b$, where the elements of δb are very small. If the matrix A has a large condition number, however, the elements of δx may not be so small. An example of this will be presented in the next section.

2.8 Hilbert matrix

We consider the Hilbert matrix H of order n , whose elements are

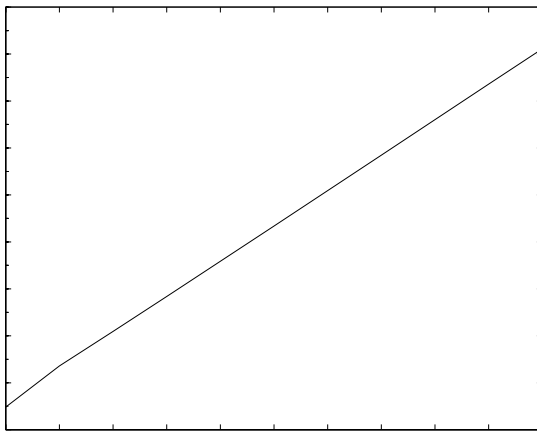
$$h_{ij} = \frac{1}{i+j-1}, \quad i, j = 1, 2, \dots, n.$$

This matrix is symmetric and positive definite (i.e., $H = H^T$, and $H^T x > 0$ for all $x \neq 0$), and therefore all of its eigenvalues are real and positive (cf. Theorem 3.1, part (ii)). However, H becomes very nearly singular as n increases. Table 2.1 shows the largest and smallest eigenvalues, and the 2-norm condition number $\kappa_2(H)$ of H , for various values of n .

n	λ_{\max}	λ_{\min}	$\kappa_2(H)$
5	1.0000	0.0001	10000
6	0.9999	0.0000	100000
7	0.9998	0.0000	1000000
8	0.9997	0.0000	10000000
9	0.9996	0.0000	100000000
10	0.9995	0.0000	1000000000
15	0.9990	0.0000	10000000000
20	0.9980	0.0000	100000000000
30	0.9960	0.0000	1000000000000
40	0.9940	0.0000	10000000000000
50	0.9920	0.0000	100000000000000
60	0.9900	0.0000	1000000000000000
70	0.9880	0.0000	10000000000000000
80	0.9860	0.0000	100000000000000000
90	0.9840	0.0000	1000000000000000000
100	0.9820	0.0000	10000000000000000000

Table 2.1. Eigenvalues and condition number of the Hilbert matrix H_n .

	max	min	κ_2^*	κ_+
(%	&&	6	5
			13	%
(&	21	%
	\$)	29	(
(\$	36	36



' 8 , $\kappa_2^* + 9$, " ! : #
& "

Figure 2.4 depicts the logarithm of the condition number $\kappa_2(H_n)$ in the 2-norm of the Hilbert matrix H_n against its order, n ; the straight line in our semilogarithmic-scale plot indicates that $\kappa_2(H_n)$, as a function of n , exhibits exponential growth. Indeed, it can be shown that

$$\kappa_2(H_n) \sim \frac{2^{n+1}}{n} \text{ as } n \rightarrow \infty.$$

We now define the vector \mathbf{b} with elements $b_j = 1/(j+1)$, $j = 1, 2, \dots, n$, chosen so that the solution of $\mathbf{A}\mathbf{x} = \mathbf{b}$, with $\mathbf{A} = H_n$, is the vector \mathbf{x} with elements $x_i = 1$, $i = 1, 2, \dots, n$. We obtain a numerical solution of the system, using the method described in Section

2.5 to give the calculated vector $\hat{x} + \delta$, and then compute the residual from $A(\hat{x} + \delta) = b + \epsilon$. The calculation uses arithmetic operations correct to 15 decimal digits, which is roughly the accuracy used by many computer systems. The results are listed in Table 2.2.

Table 2.2. Rounding errors in the solution of $Hx = b$, where H is the Hilbert matrix of order n and $b = (1, 2, \dots, n)$.

	2	2	2	2
(15	(11
)		15	&	3
(15	,	
% &		15)	
(\$	13	((2

The relative size of the residual is, in nearly every case, about the size of the basic rounding error, 10^{-15} . The resulting errors in x are smaller than the bound given by Theorem 2.11, as might be expected, since that bound corresponds to the worst possible case. In any case, for the Hilbert matrix of order greater than 14 the error is larger than the calculated solution itself, which renders the calculated solution meaningless. For matrices of this kind the condition number and the bound given by Theorem 2.11 are so large that they have little practical relevance, though they do indicate that, due to sensitivity to rounding errors, the numerical calculations are of unreliable accuracy.

The Hilbert matrix is, of course, a rather extreme example of an ill-conditioned matrix. However, we shall meet it in an important problem in Section 9.3 concerning the least squares approximation of a function by polynomials, where we shall see how a reformulation of the problem using an orthonormal basis avoids the disastrous loss of accuracy that would otherwise occur. In the next section, we introduce the idea of least squares approximation in the context of linear algebra and consider the solution of the resulting system of linear equations using the QR algorithm; this, too, relies on the notion of (ortho)normalisation.

2.9 Least squares method

Up to now, we have been dealing with systems of linear equations of the form $Ax = b$ where A is $n \times n$. However, it is frequently the case

in practical problems (typically, in problems of data-fitting) that the matrix \mathbf{A} is not square but rectangular, and we have to solve a linear system of equations $\mathbf{A}\mathbf{x} = \mathbf{b}$ with $\mathbf{A} \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$, with $m > n$; since there are more equations than unknowns, in general such a system will have no solution. Consider, for example, the linear system (with $m = 3$, $n = 2$)

$$\begin{array}{rcc} 3 & 1 & x \\ 1 & 1 & x \\ 4 & 2 & x \end{array} = \begin{array}{r} 1 \\ 0 \\ 2 \end{array};$$

by adding the first two of the three equations and comparing the result with the third, it is easily seen that there is no solution. If, on the other hand, $m < n$, then the situation is reversed and there may be an infinite number of solutions. Consider, for example, the linear system (with $m = 1$, $n = 2$)

$$(3 \ 1) \begin{array}{l} x \\ x \end{array} = 1;$$

any vector $\mathbf{x} = (\mu, 1 - 3\mu)$, with $\mu \in \mathbb{R}$, is a solution to this system.

Suppose that $m > n$; we may then need to find a vector which satisfies $\mathbf{A}\mathbf{x} = \mathbf{b}$ in \mathbb{R}^m as nearly as possible in some sense. This suggests that we define the residual vector $\mathbf{r} = \mathbf{b} - \mathbf{A}\mathbf{x}$ and require to minimise a certain norm of \mathbf{r} in \mathbb{R}^m . From the practical point of view, it is particularly convenient to minimise the residual vector \mathbf{r} in the 2-norm on \mathbb{R}^m ; this leads to the **least squares** problem:

$$\text{Minimise } \|\mathbf{r}\|_2.$$

This is clearly equivalent to minimising the square of the norm; so, on noting that

$$\|\mathbf{r}\|_2^2 = (\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}),$$

the problem may be restated as

$$\text{Minimise } (\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}).$$

Since

$$(\mathbf{b} - \mathbf{A}\mathbf{x})^T (\mathbf{b} - \mathbf{A}\mathbf{x}) = \mathbf{b}^T \mathbf{b} - 2 \mathbf{b}^T \mathbf{A}\mathbf{x} + \mathbf{x}^T \mathbf{A}^T \mathbf{A} \mathbf{x},$$

the quantity to be minimised is a nonnegative quadratic function of the n components of the vector \mathbf{x} ; the minimum therefore exists, and may

be found by equating to zero the partial derivatives with respect to the components. This leads to the system of equations

$$\mathbf{B} = \mathbf{A}^T \mathbf{A}, \quad \text{where } \mathbf{B} = \mathbf{A}^T \mathbf{A}.$$

The matrix \mathbf{B} is symmetric, and if \mathbf{A} has full rank, n , then \mathbf{B} is nonsingular; it is called the **normal** matrix, and the system $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ is called the system of **normal equations**.

The normal equations have important theoretical properties, but do not lead to a satisfactory numerical algorithm, except for fairly small problems. The difficulty is that in a practical least squares problem the matrix \mathbf{A} is likely to be quite ill-conditioned, and $\mathbf{B} = \mathbf{A}^T \mathbf{A}$ will then be extremely ill-conditioned. For example, if

$$\mathbf{A} = \begin{pmatrix} 0 \\ 0 & 1 \end{pmatrix}$$

where $\mathbf{A} = \begin{pmatrix} 0 \\ 0 & 1 \end{pmatrix}$, then $\kappa(\mathbf{A}) = \frac{1}{\epsilon} > 1$, while

$$\kappa(\mathbf{B}) = \kappa(\mathbf{A}^T \mathbf{A}) = \frac{1}{\epsilon^2} = \kappa(\mathbf{A})^2 \quad (\text{A})$$

when $0 < \epsilon < 1$. If possible, one should avoid using a method which leads to such a dramatic deterioration of the condition number.

There are various alternative techniques which avoid the direct construction of the normal matrix $\mathbf{A}^T \mathbf{A}$, and so do not lead to this extreme ill-conditioning. Here we shall describe just one algorithm, which begins by factorising the matrix \mathbf{A} , but using an orthogonal matrix rather than the lower triangular factor as in Section 2.3.

Theorem 2.12 *Suppose that $\mathbf{A} \in \mathbb{R}^{m \times n}$ where $m \geq n$. Then, \mathbf{A} can be written in the form*

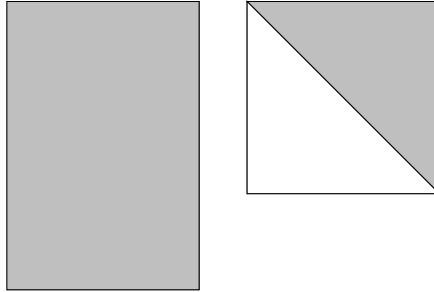
$$\mathbf{A} = \hat{\mathbf{Q}} \hat{\mathbf{R}},$$

where $\hat{\mathbf{R}}$ is an upper triangular $n \times n$ matrix, and $\hat{\mathbf{Q}}$ is an $m \times n$ matrix which satisfies

$$\hat{\mathbf{Q}}^T \hat{\mathbf{Q}} = \mathbf{I}_n, \quad (2.47)$$

where \mathbf{I}_n is the $n \times n$ identity matrix; see Figure 2.5. If $\text{rank}(\mathbf{A}) = n$, then $\hat{\mathbf{R}}$ is nonsingular.

Proof We use induction on n , the number of columns in \mathbf{A} . The theorem clearly holds when $n = 1$ so that \mathbf{A} has only one column. Indeed, writing \mathbf{a} for this column vector and assuming that $\mathbf{a} \neq \mathbf{0}$, the matrix $\hat{\mathbf{Q}}$ has just



$$A^T A = \begin{pmatrix} \sum_{i=1}^m a_{i1}^2 & \dots & \sum_{i=1}^m a_{i1} a_{i2} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m a_{i1} a_{i2} & \dots & \sum_{i=1}^m a_{i2}^2 \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^m a_{i1}^2 & \dots & \sum_{i=1}^m a_{i1} a_{i2} \\ \vdots & \ddots & \vdots \\ \sum_{i=1}^m a_{i1} a_{i2} & \dots & \sum_{i=1}^m a_{i2}^2 \end{pmatrix}$$

one column, the vector \mathbf{a}_1 , and $\hat{\mathbf{R}}$ has a single element, $\|\mathbf{a}_1\|$. In the special case where \mathbf{a}_1 is the zero vector we can choose $\hat{\mathbf{R}}$ to have the single element 0, and $\hat{\mathbf{Q}}$ to have a single column which can be an arbitrary vector in \mathbb{R}^m whose 2-norm is equal to 1.

Suppose that the theorem is true when $n = k$, where $1 \leq k < m$. Consider a matrix \mathbf{A} which has m rows and $k + 1$ columns, partitioned as

$$\mathbf{A} = (\mathbf{A} \quad \mathbf{a}_{k+1}),$$

where \mathbf{a}_{k+1} is a column vector and \mathbf{A} has k columns. To obtain the desired factorisation $\hat{\mathbf{Q}}\hat{\mathbf{R}}$ of \mathbf{A} we seek $\hat{\mathbf{Q}} = (\hat{\mathbf{Q}} \quad \mathbf{q}_{k+1})$ and

$$\hat{\mathbf{R}} = \begin{pmatrix} \hat{\mathbf{R}} \\ \mathbf{0} \end{pmatrix}$$

such that

$$\mathbf{A} = (\mathbf{A} \quad \mathbf{a}_{k+1}) = (\hat{\mathbf{Q}} \quad \mathbf{q}_{k+1}) \begin{pmatrix} \hat{\mathbf{R}} \\ \mathbf{0} \end{pmatrix}.$$

Multiplying this out and requiring that $\hat{\mathbf{Q}} \hat{\mathbf{Q}}^T = \mathbf{I}$, the identity matrix of order $k + 1$, we conclude that

$$\mathbf{A} = \hat{\mathbf{Q}} \hat{\mathbf{R}}, \tag{2.48}$$

$$= \hat{\mathbf{Q}} + \dots, \tag{2.49}$$

$$\hat{\mathbf{Q}} \hat{\mathbf{Q}}^T = \mathbf{I}, \tag{2.50}$$

$$\hat{\mathbf{Q}} = \mathbf{0}, \tag{2.51}$$

$$= 1. \tag{2.52}$$

These equations show that $\hat{Q} \hat{R}$ is the factorisation of A , which exists by the inductive hypothesis, and then lead to

$$\begin{aligned} &= \hat{Q} \quad , \\ &= (1/\alpha)(\hat{Q} \hat{Q}^T) , \end{aligned}$$

where $\alpha = \sqrt{\hat{Q} \hat{Q}^T}$. The number α is the constant required to ensure that the vector \hat{Q} is normalised.

The construction fails when $\hat{Q} \hat{Q}^T = \mathbf{0}$, for then the vector cannot be normalised. In this case we choose \hat{Q} to be any normalised vector in \mathbb{R}^n which is orthogonal in \mathbb{R}^n to all the columns of \hat{Q} , for then $\hat{Q} = \mathbf{0}$ as required. The condition at the beginning of the proof, that $k < m$, is required by the fact that when $k = m$ the matrix \hat{Q} is a square orthogonal matrix, and there is no vector \hat{Q} in \mathbb{R}^n such that $\hat{Q} = \mathbf{0}$.

With these definitions of \hat{Q} , \hat{R} , \hat{Q} and \hat{R} we have constructed the required factors of A , showing that the theorem is true when $n = k + 1$. Since it holds when $n = 1$ the induction is complete.

Now, for the final part, suppose that $\text{rank}(A) = n$. If \hat{R} were singular, there would exist a nonzero vector \hat{R} such that $\hat{R} = \mathbf{0}$; then, $A = \hat{Q} \hat{R} = \mathbf{0}$, and hence $\text{rank}(A) < n$, contradicting our hypothesis that $\text{rank}(A) = n$. Therefore, if $\text{rank}(A) = n$, then \hat{R} is nonsingular. □

The matrix factorisation whose existence is asserted in Theorem 2.12 is called the **QR factorisation**. Here, we shall present its use in the solution of least squares problems. In Chapter 5 we shall revisit the idea in a different context which concerns the numerical solution of eigenvalue problems.

Theorem 2.13 *Suppose that $A \in \mathbb{R}^{m \times n}$, with $m \geq n$ and $\text{rank}(A) = n$, and let $\hat{Q} \in \mathbb{R}^{m \times n}$. Then, there exists a unique least squares solution of the system of equations $Ax = b$: a vector x in \mathbb{R}^n which minimises the function $\|Ax - b\|$ over all x in \mathbb{R}^n . The vector x can be obtained by finding the factors \hat{Q} and \hat{R} of A defined in Theorem 2.12, and then solving the nonsingular upper triangular system $\hat{R}x = \hat{Q}^T b$.*

Proof The matrix \hat{Q} has m rows and n columns, with $m \geq n$, and it satisfies

$$\hat{Q} \hat{Q}^T = I_n .$$

We shall suppose that $m > n$, the case $m = n$ being a trivial special case with

$$= A^{-} = (\hat{Q}\hat{R})^{-} = \hat{R}^{-} \hat{Q}^{-} = \hat{R}^{-} \hat{Q}^{-} ,$$

and hence $\hat{R}^{-} = \hat{Q}^{-}$, as required.

For $m > n$ now, the vector $\hat{R}^{-} \hat{Q}^{-}$ can be written as the sum of two vectors:

$$= \hat{R}^{-} \hat{Q}^{-} + \hat{R}^{-} \hat{Q}^{-} ,$$

where $\hat{R}^{-} \hat{Q}^{-}$ is in the linear space spanned by the n columns of the matrix \hat{Q} , and $\hat{R}^{-} \hat{Q}^{-}$ is in the orthogonal complement of this space in \mathbb{R}^m . The vector $\hat{R}^{-} \hat{Q}^{-}$ is a linear combination of the columns of \hat{Q} , and $\hat{R}^{-} \hat{Q}^{-}$ is orthogonal to every column of \hat{Q} ; *i.e.*, there exists $\hat{R}^{-} \hat{Q}^{-}$ such that

$$= \hat{R}^{-} \hat{Q}^{-} + \hat{R}^{-} \hat{Q}^{-} , \quad \hat{R}^{-} \hat{Q}^{-} = \hat{Q}^{-} , \quad \hat{Q}^{-} \hat{Q}^{-} = \mathbf{0} . \tag{2.53}$$

Now, suppose that $\hat{R}^{-} \hat{Q}^{-}$ is the solution of $\hat{R}^{-} \hat{Q}^{-} = \hat{Q}^{-}$, and that $\hat{R}^{-} \hat{Q}^{-}$ is any vector in \mathbb{R}^m . Then,

$$\begin{aligned} A \hat{R}^{-} \hat{Q}^{-} &= \hat{Q} \hat{R}^{-} \hat{Q}^{-} \\ &= \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + \hat{Q} \hat{R}^{-} \hat{Q}^{-} \\ &= \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + \hat{Q} \hat{Q}^{-} \\ &= \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + \hat{Q} \hat{Q}^{-} + \hat{Q} \hat{Q}^{-} \\ &= \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + \hat{Q} \hat{Q}^{-} \hat{Q}^{-} \\ &= \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) , \end{aligned}$$

where we have used (2.53) repeatedly; in particular, the last equality follows by noting that $\hat{Q}^{-} \hat{Q}^{-} = \mathbf{I}$. Hence

$$\begin{aligned} A \hat{R}^{-} \hat{Q}^{-} &= (\hat{R}^{-} \hat{Q}^{-}) \hat{R}^{-} \hat{Q}^{-} \hat{Q} \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + 2(\hat{R}^{-} \hat{Q}^{-}) \hat{R}^{-} \hat{Q}^{-} \\ &= \hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) + \end{aligned}$$

since $\hat{Q}^{-} \hat{Q}^{-} = \mathbf{0}$. Thus $A \hat{R}^{-} \hat{Q}^{-}$ is smallest when $\hat{R}^{-} (\hat{R}^{-} \hat{Q}^{-}) = \mathbf{0}$, which implies that $\hat{R}^{-} \hat{Q}^{-} = \mathbf{0}$, since the matrix \hat{R}^{-} is nonsingular. Hence $\hat{R}^{-} \hat{Q}^{-}$, defined as the solution of $\hat{R}^{-} \hat{Q}^{-} = \hat{Q}^{-}$, is the required least squares solution. \square

2.10 Notes

There are many good books on the subject of numerical linear algebra which cover the topics discussed in this chapter in much greater detail,

and address questions which we have not touched on here. Without any attempt to be exhaustive, we single out four texts from the vast literature. The first two books on the list below are well-known monographs on the subject, while the last two are excellent textbooks.

 , *Matrix Computations*, Third Edition, Johns Hopkins University Press, Baltimore, 1996.

 , *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.

 ! "###, *Numerical Linear Algebra*, SIAM, Philadelphia, 1997.

\$, *Introduction to Numerical Linear Algebra and Optimisation*, Cambridge University Press, Cambridge, 1989.

As we have already noted in Section 2.2, the invention of the elimination technique is attributed to Gauss who published the method in his *Theoria motus* (1809), although the idea was already known to the Chinese two thousand years ago. Gauss himself was concerned with positive definite systems. The method was extended to linear systems with general matrices by Jacobi. The interpretation of Gaussian elimination as matrix factorisation is due to P.S. Dwyer: A matrix presentation of least squares and correlation theory with matrix justification of improved methods of solutions, *Ann. Math. Stat.* **15**, 82–89, 1944.

The sensitivity of Gaussian elimination to rounding errors was studied by Wilkinson in Error analysis of direct methods of matrix inversion, *J. Assoc. Comput. Math.* **8**, 281–330, 1961. The idea of pivoting was used as early as 1947 by von Neumann and Goldstein. The concept of the condition number of a matrix was introduced by Turing in Rounding-off errors in matrix processes, *Quart. J. Mech. Appl. Math.* **1**, 287–308, 1948. Our treatment of condition numbers follows the textbook of Trefethen and Bau, cited above.

1 F 5K, K, B# /, #. O 6 \$
B F 6D C #. 9 6 #.># F 6D G 4 5 \$
, 6 F4 4 1K, || F , 5
1)
2 K O 6 " B & ##># C > H,
#@. ? : D4
3 K 5 B . /, #@ * % CO 6 B
O 6D C . 9 6 #@>& / %D4
4 6 6 F B* /, #@ * O C K 6 #@. @
, %D4
5 % , B * K ##># : C & K #@>
D4

Normed linear spaces play a key role in functional analysis (see, for example, K. Yosida, *Functional Analysis*, Third Edition, Springer, Berlin, 1971, page 30). Here, we have concentrated on finite-dimensional normed linear spaces over the field of real numbers.

The relevance of norms in numerical linear algebra was highlighted by Householder in his book *The Theory of Matrices in Numerical Analysis*, Blaisdell, New York, 1964.

The idea of least squares fitting is due to Gauss, who invented the method in the 1790s. However, it was the French mathematician Legendre who first published the method in 1806 in a book on determining the orbits of comets. Legendre’s method involved a number of observations taken at equal intervals and he assumed that the comet followed a parabolic path, so he ended up with more equations than there were unknowns. Legendre then applied his methods to the data known for two comets. In an Appendix to the book Legendre described the least squares method of fitting a curve to the data available. Gauss published his version of the least squares method in 1809 and, although acknowledging that it had already appeared in Legendre’s book, Gauss nevertheless claimed priority for himself. This greatly hurt Legendre, leading to one of the infamous priority disputes in the history of mathematics. A recent exhaustive monograph on numerical algorithms for least squares problems is due to Å. Björk: *Numerical Methods for Least Squares Problems*, SIAM, Philadelphia, 1996.

The version of the QR factorisation considered here is the *reduced version*, following the terminology in Chapter 7 of Trefethen and Bau. In the *full version* of the QR factorisation for a matrix $A \in \mathbb{R}^{m \times n}$, we have $A = QR$, where $Q \in \mathbb{R}^{m \times m}$, $R \in \mathbb{R}^{m \times n}$ (cf. Chapter 5).

In a footnote to Definition 2.12 we mentioned the Moore–Penrose generalised inverse A^+ of a matrix $A \in \mathbb{R}^{m \times n}$. A^+ can be defined through the singular value decomposition of A (cf. L.N. Trefethen and D. Bau, III: *Numerical Linear Algebra*, SIAM, Philadelphia, 1997). Recall that the singular values of A are the square roots of the (nonnegative) eigenvalues of the matrix $A^T A$.

¹ % , O B> ' 6 #@ , "1) % C K 6 #@@*
 ' 1 %D 1 1 , ' (+ 5 1 4
 O || 6 6 F4 4 , ' (+ 5 1 4

² % \$: B# #&> 9 , C # K 6 #. **
 9 , D4

Theorem 2.14 (Singular value decomposition) Let $A \in \mathbb{R}^{m \times n}$; then, there exist $U \in \mathbb{R}^{m \times m}$, $\Sigma \in \mathbb{R}^{m \times n}$ and $V \in \mathbb{R}^{n \times n}$ such that

$$A = U\Sigma V^T,$$

where Σ is a diagonal matrix whose diagonal entries, σ_i , $i = 1, 2, \dots, n$, are the singular values of A , $U^T U = I$ and $V^T V = I$, with I denoting the $n \times n$ identity matrix.

The Moore–Penrose generalised inverse of the diagonal matrix $\Sigma \in \mathbb{R}^{m \times n}$ is defined as the diagonal matrix $\Sigma^+ \in \mathbb{R}^{n \times m}$ whose diagonal entries are

$$\sigma_i^{-1} \text{ if } \sigma_i > 0, \\ 0 \text{ if } \sigma_i = 0.$$

The generalised inverse $A^+ \in \mathbb{R}^{n \times m}$ of a matrix $A \in \mathbb{R}^{m \times n}$ with singular value decomposition $A = U\Sigma V^T$ is defined by

$$A^+ = V\Sigma^+ U^T.$$

In the special case when $m = n$ and $A \in \mathbb{R}^{n \times n}$ is nonsingular, the n singular values of A are all nonzero and therefore $\Sigma^+ = \Sigma^{-1}$. Hence, also, $A^+ = A^{-1}$, which then justifies the use of the terminology ‘generalised inverse’ for the matrix A^+ defined above.

Exercises

2.1 Let $n \geq 2$. Given the matrix $A = (a_{ij}) \in \mathbb{R}^{n \times n}$, the permutation matrix $Q \in \mathbb{R}^{n \times n}$ reverses the order of the rows of A , so that $(QA)_{ij} = a_{(n-i+1)j}$. If $L \in \mathbb{R}^{n \times n}$ is a lower triangular matrix, what is the structure of the matrix QLQ ?

Show how to factorise $A \in \mathbb{R}^{n \times n}$ in the form $A = UL$, where $U \in \mathbb{R}^{n \times n}$ is unit upper triangular and $L \in \mathbb{R}^{n \times n}$ is lower triangular. What conditions on A will ensure that the factorisation exists? Give an example of a square matrix A which cannot be factorised in this way.

2.2 Let $n \geq 2$. Consider a matrix $A \in \mathbb{R}^{n \times n}$ whose every leading principal submatrix of order less than n is nonsingular. Show that A can be factored in the form $A = LDU$, where $L \in \mathbb{R}^{n \times n}$ is unit lower triangular, $D \in \mathbb{R}^{n \times n}$ is diagonal and $U \in \mathbb{R}^{n \times n}$ is unit upper triangular.

If the factorisation $A = LU$ is known, where L is unit lower

triangular and U is upper triangular, show how to find the factors of the transpose A^T .

- 2.3 Let $n \geq 2$ and suppose that the matrix $A \in \mathbb{R}^{n \times n}$ is nonsingular. Show by induction, as in Theorem 2.3, that there are a permutation matrix $P \in \mathbb{R}^{n \times n}$, a lower triangular matrix $L \in \mathbb{R}^{n \times n}$, and a unit upper triangular matrix $U \in \mathbb{R}^{n \times n}$ such that $PA = LU$.

By finding a suitable 2×2 matrix A , or otherwise, show that this may not be true if A is singular.

- 2.4 The lower triangular matrix $L \in \mathbb{R}^{n \times n}$, $n \geq 2$, is nonsingular, and the vector $b \in \mathbb{R}^n$ is such that $b_i = 0$, $i = 1, 2, \dots, k$, with $1 \leq k < n$. The vector $y \in \mathbb{R}^n$ is the solution of $Ly = b$. Show, by partitioning L , that $y_j = 0$, $j = 1, 2, \dots, k$. Hence give an alternative proof of Theorem 2.1(iv), that the inverse of a nonsingular lower triangular matrix is itself lower triangular.

- 2.5 Given a matrix $A \in \mathbb{R}^{n \times n}$, define the matrix $B \in \mathbb{R}^{n \times n}$ in which the first n columns are the columns of A , and the last n columns are the columns of the identity matrix I . Consider the following computational scheme. Treat the rows of the matrix B in order, so that $j = 1, 2, \dots, n$. Multiply every element in row j by the reciprocal of the diagonal element, $1/b_{jj}$; then, replace every element b_{ij} which is not in row j , so that $i \neq j$, by $b_{ij} - b_{ij}b_{jj}^{-1}b_{jj}$.

Show that the result is equivalent to multiplying B on the left by a sequence of matrices. Explain why, at the end of the computation, the first n columns of B are the columns of the identity matrix I , and the last n columns are the columns of the inverse matrix A^{-1} . Give a condition on the matrix A which will ensure that the computation does not break down.

Show that the process as described requires approximately $2n^2$ multiplications, but that, if the multiplications in which one of the factors is zero are not counted, the total is approximately n^3 .

- 2.6 Use the method of Exercise 5 to find the inverse of the matrix

$$A = \begin{pmatrix} 2 & 4 & 2 \\ 1 & 0 & 3 \\ 3 & 1 & 2 \end{pmatrix}.$$

2.7 Suppose that for a matrix $A \in \mathbb{R}^{n \times n}$,

$$|a_{jj}| \geq \sum_{i \neq j} |a_{ij}|, \quad j = 1, 2, \dots, n.$$

Show that, for any vector $x \in \mathbb{R}^n$,

$$\|(Ax)\|_1 \leq \|x\|_1.$$

Find a nonzero vector x for which equality can be achieved, and deduce that

$$\|A\|_1 = \max_{1 \leq j \leq n} \sum_{i=1}^n |a_{ij}|.$$

2.8 (i) Show that, for any vector $x = (x_1, \dots, x_n)^T$,

$$\|x\|_1 \leq \|x\|_2 \leq \sqrt{2} \|x\|_1.$$

In each case give an example of a nonzero vector x for which equality is attained. Deduce that $\|x\|_1 \leq \|x\|_2 \leq \sqrt{2} \|x\|_1$. Show also that $\|x\|_2 \leq \sqrt{2} \|x\|_1$.

(ii) Show that, for any matrix $A \in \mathbb{R}^{n \times n}$,

$$\|A\|_1 \geq \bar{\sigma}(A) \quad \text{and} \quad \|A\|_2 \geq \bar{\sigma}(A).$$

In each case give an example of a matrix A for which equality is attained. (See the footnote following Definition 2.12 for the meaning of $\bar{\sigma}(A)$, $\sigma(A)$ and $\|A\|_2$ when $A \in \mathbb{R}^{n \times n}$.)

2.9 Prove that, for any nonsingular matrix $A \in \mathbb{R}^{n \times n}$,

$$\kappa(A) = \frac{\|A\|_2}{\sigma_{\min}(A)},$$

where $\sigma_{\min}(A)$ is the smallest and $\sigma_{\max}(A)$ is the largest eigenvalue of the matrix $A^T A$.

Show that the condition number $\kappa(Q)$ of an orthogonal matrix Q is equal to 1. Conversely, if $\kappa(A) = 1$ for the matrix A , show that all the eigenvalues of $A^T A$ are equal; deduce that A is a scalar multiple of an orthogonal matrix.

2.10 Let $A \in \mathbb{R}^{n \times n}$. Show that if λ is an eigenvalue of $A^T A$, then

$$0 \leq \lambda \leq \|A\|_2^2,$$

provided that the same subordinate matrix norm is used for

both A and A^{-1} . Hence show that, for any nonsingular $n \times n$ matrix A ,

$$(A^{-1})^{-1} = A \quad (A^{-1})^{-1} = A^{-1}.$$

- 2.11 For the matrix defined by (2.46) write down the matrix $A^{-1}A$. Show that any vector $x = 0$ is an eigenvector of $A^{-1}A$ with eigenvalue $\lambda = 1$, provided that $x_1 = 0$ and $x_2 + x_3 + x_4 = 0$. Show also that there are two eigenvectors with $x_1 = x_2 = x_3$ and find the corresponding eigenvalues. Deduce that

$$(A^{-1}A)^{-1} = -(n+1) \begin{bmatrix} 1 & & & \\ & 1 & & \\ & & 1 & \\ & & & 1 \end{bmatrix}.$$

- 2.12 Let $B \in \mathbb{R}^{n \times n}$ and denote by I the identity matrix of order n . Show that if the matrix $I - B$ is singular, then there exists a nonzero vector x such that $(I - B)x = 0$; deduce that $\|B\| \geq 1$, and hence that, if $\|A\| < 1$, then the matrix $I - A$ is nonsingular.

Now suppose that $A \in \mathbb{R}^{n \times n}$ with $\|A\| < 1$. Show that

$$(I - A)^{-1} = I + A(I - A)^{-1},$$

and hence that

$$(I - A)^{-1} = \frac{1}{1 - \|A\|} (I - A)^{-1}.$$

Deduce that

$$(I - A)^{-1} = \frac{1}{1 - \|A\|}.$$

- 2.13 Let $A \in \mathbb{R}^{n \times n}$ be a nonsingular matrix and $\|A\| < 1$. Suppose that $\|A^{-1}\| = \frac{1}{1 - \|A\|}$ and $(A + A^{-1})^{-1} = \frac{1}{2}$, and that $\|A^{-1}\| - \|A\| < 1$. Use the result of Exercise 12 to show that

$$\frac{1}{\|A^{-1}\|} = \frac{\|A^{-1}\| - \|A\|}{1 - \|A\|}.$$

- 2.14 Suppose that $A \in \mathbb{R}^{n \times n}$ is a nonsingular matrix, and $\|A\| < 1$. Given that $\|A^{-1}\| = \frac{1}{1 - \|A\|}$ and $\|A^{-1}\| - \|A\| = \frac{1}{2}$, Theorem 2.11 states that

$$\|A\| = \frac{1}{2}.$$

By considering the eigenvectors of $A^{-1}A$, show how to find vectors x and y for which equality is attained, when using the 2-norm.

2.15 Find the QR factorisation of the matrix

$$A = \begin{pmatrix} 9 & 6 \\ 12 & 8 \\ 0 & 20 \end{pmatrix},$$

and hence find the least squares solution of the system of linear equations

$$\begin{aligned} 9x + 6y &= 300, \\ 12x + 8y &= 600, \\ 20y &= 900. \end{aligned}$$

3

Special matrices

3.1 Introduction

In this chapter we show how one can modify the elimination method for the solution of $Ax = b$ when the matrix A has certain special properties. In particular when A is symmetric and positive definite the amount of computational work can be halved. For matrices with a band structure, having nonzero elements only in positions close to the diagonal, the efficiency can be improved even more dramatically.

3.2 Symmetric positive definite matrices

Definition 3.1 The matrix $A = (a_{ij})_{n \times n}$ is said to be **symmetric** if $a_{ij} = a_{ji}$ for all i and j in the set $\{1, 2, \dots, n\}$; i.e., if $A = A^T$. The set of all symmetric matrices A will be denoted by S_n . A matrix A is called **positive definite** if

$$x^T A x > 0$$

for every vector $x \neq 0$.

Example 3.1 Consider the matrix $A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$,

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

and a vector $x = (x_1, x_2)^T \neq 0$.

Clearly, $x^T A x = ax_1^2 + (b+c)x_1x_2 + dx_2^2$. The quadratic form on the right-hand side is positive for all real numbers x_1, x_2 such that

$$= (x, x) = (0, 0) = \mathbf{0} \text{ if, and only if,}$$

$$a > 0, \quad d > 0 \text{ and } (b + c) < 4ad.$$

We see that if A is positive definite, then the diagonal elements of A are positive. Further, noting that the third inequality can be rewritten as

$$(b - c) < 4(ad - bc) = 4 \det(A),$$

we deduce that the determinant of a positive definite matrix A is positive. This, of course, is still true in the special case when A is symmetric, *i.e.*, when $b = c$.

The next theorem extends the observations of the last example to any symmetric positive definite matrix A .

Theorem 3.1 *Suppose that $n \geq 2$ and $A = (a_{ij})$ is positive definite; then:*

- (i) *all the diagonal elements of A are positive, that is, $a_{ii} > 0$, for $i = 1, 2, \dots, n$;*
- (ii) *all the eigenvalues of A are real and positive, and the eigenvectors of A belong to \mathbb{R}^n ;*
- (iii) *the determinant of A is positive;*
- (iv) *every submatrix B of A obtained by deleting any set of rows and the corresponding set of columns from A is symmetric and positive definite; in particular, every leading principal submatrix is positive definite;*
- (v) *$a_{ii} < a_{jj}$ for all i and j in $1, 2, \dots, n$ such that $i < j$;*
- (vi) *the element of A with largest absolute value lies on the diagonal;*
- (vii) *if a_{ii} is the largest of the diagonal elements of A , then*

$$a_{ii} \geq a_{jj} \quad i, j = 1, 2, \dots, n.$$

Proof (i) Consider the vector x with only one nonzero element, in position $i = 1, 2, \dots, n$. Since A is positive definite and $x \neq 0$, it follows that $x^T A x = a_{ii} > 0$, and therefore $a_{ii} > 0$.

(ii) Suppose that λ is an eigenvalue of A and let $x = \mathbf{0}$ denote the associated eigenvector. Further, let \bar{x} denote the vector in \mathbb{R}^n whose i th element is the complex conjugate of the i th element of

, $i = 1, 2, \dots, n$. As $\mathbf{A}^{-1} = (\mathbf{A}^{-1})^T$, it follows that $\lambda^{-1} \mathbf{A}^{-1} = (\lambda^{-1} \mathbf{A}^{-1})^T$, and therefore, using the symmetry of \mathbf{A} ,

$$\mathbf{A}^{-1} = \mathbf{A}^{-1} = (\lambda^{-1} \mathbf{A}^{-1}) = (\lambda^{-1} \mathbf{A}^{-1})^T = (\lambda^{-1} \mathbf{A}^{-1})^T.$$

Complex conjugation then yields $\overline{\lambda^{-1} \mathbf{A}^{-1}} = \overline{\lambda^{-1} \mathbf{A}^{-1}}$, and hence $\overline{\lambda^{-1} \mathbf{A}^{-1}} = \overline{\lambda^{-1} \mathbf{A}^{-1}}$. As $\mathbf{A}^{-1} = \mathbf{0}$, it follows that $\lambda = \overline{\lambda}$; i.e., λ is a real number.

The fact that the eigenvector associated with λ has real elements follows by noting that all elements of the singular matrix $\mathbf{A} - \lambda \mathbf{I}$ are real numbers. Therefore, the column vectors of $\mathbf{A} - \lambda \mathbf{I}$ are linearly dependent in \mathbb{R}^n . Hence there exist n real numbers x_1, \dots, x_n such that $(\mathbf{A} - \lambda \mathbf{I}) \mathbf{x} = \mathbf{0}$, where $\mathbf{x} = (x_1, \dots, x_n)^T$.

Finally, as $\mathbf{A} = \mathbf{A}^T$ with λ real and \mathbf{x} real, we have that $\mathbf{A} \mathbf{x} = \lambda \mathbf{x}$. Since $\lambda = \mathbf{A} / \mathbf{x}^T \mathbf{x}$ and \mathbf{A} is positive definite, λ is the ratio of two positive real numbers and therefore also real and positive.

(iii) This follows from the fact that the determinant of \mathbf{A} is equal to the product of its eigenvalues, and the previous result. Indeed, since \mathbf{A} is symmetric, there exist an orthogonal matrix \mathbf{X} and a diagonal matrix Λ , whose diagonal elements are the eigenvalues $\lambda_i, i = 1, 2, \dots, n$, of \mathbf{A} , such that $\mathbf{A} = \mathbf{X} \Lambda \mathbf{X}^T = \mathbf{X}^{-1} \Lambda \mathbf{X}$. By the Binet–Cauchy Theorem (see Chapter 2, end of Section 2.3),

$$\begin{aligned} \det(\mathbf{A}) &= \det(\mathbf{X}^{-1}) \det(\Lambda) \det(\mathbf{X}) \\ &= \frac{1}{\det(\mathbf{X})} \det(\Lambda) \det(\mathbf{X}) \\ &= \det(\Lambda) = \lambda_1 \dots \lambda_n > 0. \end{aligned}$$

(iv) Consider the vector \mathbf{v} with zeros in the positions corresponding to the rows which have been deleted. Then,

$$\mathbf{A} \mathbf{v} = \mathbf{B} \mathbf{w}$$

where \mathbf{B} is the submatrix of \mathbf{A} containing the rows and columns which remain after deletion, and \mathbf{w} is the vector consisting of the elements of \mathbf{v} which were not deleted. Since the expression on the left is positive, the same is true of the expression on the right, for all vectors \mathbf{w} except the zero vector. Therefore \mathbf{B} is positive definite.

(v) By the previous result the 2×2 submatrix consisting of rows and columns r and s of \mathbf{A} is positive definite, and its determinant is therefore positive.

(vi) This follows from the previous result, since it shows that a_{rr} cannot exceed the greater of a_{ss} and a_{tt} .

(vii) This follows at once from the previous result. □

The converses of two of these results are also true:

- (i) If all the eigenvalues of the symmetric matrix A are positive, then A is positive definite;
- (ii) If the determinant of each leading principal submatrix of a matrix A is positive, then A is positive definite.

The proof of the second result is involved and will not be given here; see, however, Example 3.1 for the case of $n = 2$. The proof of the first statement, on the other hand, is quite simple and proceeds as follows.

Since A is symmetric, it has a complete set of orthonormal eigenvectors v_1, \dots, v_n in \mathbb{R}^n , and the corresponding eigenvalues $\lambda_1, \dots, \lambda_n$ are all real. Given any vector $x \in \mathbb{R}^n$, it can be expressed as

$$x = \sum_{i=1}^n \alpha_i v_i$$

where $\alpha_i = v_i^T x$, $i = 1, 2, \dots, n$, and $\alpha_i^2 \geq 0$. Since $A v_i = \lambda_i v_i$, $i = 1, 2, \dots, n$, it follows that

$$A x = \sum_{i=1}^n \alpha_i \lambda_i v_i$$

As $\alpha_i = 0$ for $i = j$ and $\alpha_i^2 = 1$, we deduce that

$$A x = \sum_{i=1}^n \lambda_i \alpha_i^2 v_i$$

$$\min_{i=1, \dots, n} \lambda_i > 0,$$

since $\min_{i=1, \dots, n} \lambda_i > 0$; therefore A is positive definite.

For a symmetric positive definite matrix A we can now obtain an LU factorisation $A = LU$ in which $U = L^T$.

Theorem 3.2 Suppose that $n \geq 2$ and A is a positive definite matrix; then, there exists a lower triangular matrix L such that

$$A = LL^T.$$

This is known as the **Cholesky factorisation** of A .

1 9 4%4 O 4 4 K ' %
 5 6 #@@ &4 4>4
 2 P% -\$: "6 B#. &>C#@#. D 9 , 6 S, 5 5
 6 5 6 %1 , 7 1)4 O

Proof Since A is symmetric and positive definite, all the leading principal submatrices of A are positive definite, and hence by Theorem 2.2 the usual LU factorisation exists, with

$$A = L U,$$

L a unit lower triangular and U an upper triangular matrix. In this factorisation the product of the leading principal submatrices of L and U of order k is the leading principal submatrix of A of order k , $1 \leq k \leq n$. Since the determinant of this submatrix is positive and all the diagonal elements of L are unity, it follows that

$$u_{11} u_{22} \dots u_{kk} > 0, \quad k = 1, 2, \dots, n.$$

Thus all the diagonal elements of U are positive. If we now define D to be the diagonal matrix with elements $d_{ii} = u_{ii}$, $i = 1, 2, \dots, n$, we can write

$$A = L U = (L D)(D^{-1} U) = L U,$$

where now $l_{ii} = u_{ii} = d_{ii}$. The symmetry of the matrix A shows that

$$LU = A = A^T = U^T L^T,$$

so that

$$U(L^{-1})^T = L^{-1} U^T.$$

In this equality the left-hand side is upper triangular, and the right-hand side is lower triangular, and hence both sides must be diagonal. Therefore, $U = D L^{-1}$, where D is a diagonal matrix; but U and L have the same diagonal elements, so $D = I$ and $U = L^{-1}$.

The same argument shows that L and L^{-1} are unique, except for the arbitrary choice of the signs of the square roots in the definition of the diagonal matrix D . If we make the natural choice, taking all the square roots to be positive, then the diagonal elements of L are positive, and the factorisation is unique. □

5
G 1 G 1
6
23 3 / 4 C 5 ' (1) 9 6 #. #@@ (S, 6 O "

In practice we construct the elements of L directly, rather than forming L^{-1} and U first. This is done in a similar way to the LU factorisation. Suppose that $i < j$; we then require that

$$a_{ij} = l_{ij} l_{jj}, \quad 1 \leq i < j \leq n. \quad (3.1)$$

Note that we have used the fact that $(L^{-1})_{ij} = l_{ij}$; the sum only extends up to $k = i$ since L is lower triangular. The same equation will also hold for $i > j$, since A is symmetric. For $i = j$, equation (3.1) gives

$$l_{ii} = a_{ii}^{-1/2}, \quad l_{ij} = a_{ij} a_{ii}^{-1/2}, \quad 1 < i \leq j \leq n. \quad (3.2)$$

As A is a positive definite matrix, $a_{ii} > 0$ and therefore l_{ii} is a positive real number. Further, as we have seen in the proof of the preceding theorem, $l_{ij} > 0$, $i = 2, 3, \dots, n$. We find similarly that

$$l_{ij} = \frac{1}{l_{jj}} a_{ij} - l_{ik} l_{kj}, \quad 1 \leq i < j \leq n. \quad (3.3)$$

These equations now enable us to calculate the elements of L in succession. For each $i = 1, 2, \dots, n-1$, we first calculate l_{ii} from (3.2), and then calculate $l_{i2}, l_{i3}, \dots, l_{in}$ from (3.3). Finally, we compute l_{ij} using (3.2).

As, by hypothesis, the matrix A^{-1} is positive definite, the required factorisation exists, so we can be sure that the divisor l_{ii} in (3.3), and the expression in the curly brackets in (3.2) whose square root is taken, will be positive. Thus, (3.2) implies that

$$l_{ij} = a_{ij} / \max\{l_{ik} l_{kj}, a_{ij}\}, \quad i = 2, 3, \dots, n.$$

The elements of the factor L cannot therefore grow very large, and no pivoting is necessary.

The evaluation of l_{ij} from (3.2) requires $i-1$ multiplications, $i-1$ subtractions and one square root operation, a total of $2i-1$ operations. The calculation of each l_{ij} from (3.3) also requires $2i-1$ operations. The total number of operations required to construct L is therefore

$$(2i-1) = (2i-1)(1+n-i) = -n(n+1)(2n+1).$$

For large n the number of operations required is approximately $\frac{1}{2}n^2$, which, as might be expected, is half the number given in Section 2.6 for the LU factorisation of a nonsymmetric matrix.

3.3 Tridiagonal and band matrices

As we shall see in the final chapters, in the numerical solution of boundary value problems for second-order differential equations one encounters a particular kind of matrix whose elements are mostly zeros, except for those along its main diagonal and the two adjacent diagonals. Matrices of this kind are referred to as tridiagonal. In order to motivate the definition of tridiagonal matrix stated in Definition 3.2 below, we begin with an example which is discussed in more detail in Chapter 13.

Example 3.2 Consider the two-point boundary value problem

$$\frac{d y}{d x} + r(x)y = f(x), \quad x \in (0, 1),$$

$$y(0) = 0, \quad y(1) = 0.$$

where r and f are continuous functions of x defined on the interval $[0, 1]$.

The numerical solution of the boundary value problem proceeds by selecting an integer $n \geq 4$, choosing a step size $h = 1/n$, and subdividing the interval $[0, 1]$ by the points $x_k = kh$, $k = 0, 1, \dots, n$. The numerical approximation to $y(x_k)$, the value of the analytical solution y at the point $x = x_k$, is denoted by Y_k . The values Y_k are obtained by solving the set of linear equations

$$\frac{Y_{k+1} - 2Y_k + Y_{k-1}}{h} + r(x_k)Y_k = f(x_k)$$

for $k = 1, 2, \dots, n-1$, together with the boundary conditions

$$Y_0 = 0, \quad Y_n = 0.$$

Equivalently,

$$\begin{aligned} a Y_{k-1} + c Y_k + b Y_{k+1} &= d, \quad k = 1, 2, \dots, n-1, \\ Y_0 &= 0, \quad Y_n = 0, \end{aligned}$$

where

$$a_k = b_k = 1/h, \quad c_k = 2/h + r(x_k), \quad d_k = f(x_k),$$

for $k = 1, 2, \dots, n - 1$.

Clearly, for $1 < k < n - 1$, the k th equation in the linear system above involves only three of the $n - 1$ unknowns: Y_{k-1} , Y_k and Y_{k+1} .

The example motivates the following definition of a tridiagonal (or triple diagonal) matrix.

Definition 3.2 *Suppose that $n \geq 3$. A matrix $T = (t_{ij})_{n \times n}$ is said to be **tridiagonal** if it has nonzero elements only on the main diagonal and the two adjacent diagonals; i.e.,*

$$t_{ij} = 0 \quad \text{if} \quad |i - j| > 1, \quad i, j = 1, 2, \dots, n.$$

Such matrices are also sometimes called **triple diagonal**.

It is easy to see that in the LU factorisation process of a tridiagonal matrix $T_{n \times n}$, without row interchanges, the unit lower triangular matrix $L_{n \times n}$ and the upper triangular matrix $U_{n \times n}$ each have only two elements in each row. Writing T in the compact notation

$$T = \begin{pmatrix} & b & & c & & & \\ a & & b & & c & & \\ & a & & b & & c & \\ \dots & \dots & \dots & \dots & \dots & \dots & \dots \\ \dots & \dots & \dots & \dots & \dots & a & b \end{pmatrix}, \quad (3.4)$$

the factorisation may be written $T = LU$ where

$$L = \begin{pmatrix} & & & & & & \\ & 1 & & & & & \\ & & l & & & & \\ & & & 1 & & & \\ \dots & \dots & \dots & \dots & \dots & \dots & \\ & & & & & l & \\ & & & & & & 1 \end{pmatrix} \quad (3.5)$$

and

$$U = \begin{pmatrix} u & & & & & & \\ & v & & & & & \\ & & u & & & & \\ & & & v & & & \\ & & & & u & & v \\ \dots & \dots & \dots & \dots & \dots & \dots & \\ & & & & & & u \end{pmatrix}, \quad (3.6)$$

with the convention that the missing elements in these matrices are all equal to zero. It is often convenient to define $a_n = 0$ and $c_1 = 0$. Multiplying L and U shows that $v_j = c_j$, and that the elements l_j and u_j can be calculated from

$$l_j = a_j / u_{j-1}, \quad u_j = b_j - l_j c_{j-1}, \quad j = 2, 3, \dots, n, \quad (3.7)$$

starting from $u_1 = b_1$.

Let us suppose that our aim is to solve the system of linear equations $Tx = r$, where the matrix T is tridiagonal and nonsingular, and r is a vector. Having calculated the elements of the matrices L and U in the LU factorisation $T = LU$ using (3.7), the forward and backsubstitution are then also very simple. Letting $y = Ux$, the equation $Lx = r$ gives

$$y_j = r_j, \quad (3.8)$$

$$y_j = r_j - l_j y_{j-1}, \quad j = 2, 3, \dots, n, \quad (3.9)$$

and finally from $Ux = y$ we get

$$x_j = y_j / u_j, \quad (3.10)$$

$$x_j = (y_j - v_j x_{j+1}) / u_j, \quad j = n-1, n-2, \dots, 1. \quad (3.11)$$

The LU factorisation of a tridiagonal matrix requires approximately $3n$ operations. The forward and backsubstitution together involve approximately $5n$ operations. Thus, the whole solution process requires approximately $8n$ operations. The total amount of work is therefore far less than for a full matrix, being of order n^3 for large n , compared with n^3 for a full matrix. The method we have described is a minor variation on what is often known as the *Thomas algorithm*.

So far we have assumed that pivoting was not necessary; clearly any interchange of rows will destroy the tridiagonal structure of T . However, it is easy to see that the only interchanges required will be between two adjacent rows.

Theorem 3.3 Suppose that $n \geq 3$ and T is a tridiagonal matrix; then, there exists a permutation matrix P such that

$$PA = LU \quad (3.12)$$

1 %1 : 6 O4 5 6) ' || 6 , #> \$
 6 " 6 P ' , C9 , 64 O 4
 . 4 : 4) # 5 ./ 1 / 5 ! 4 : 4O4
 @ 4 : 4 5 6 ! " # @ 4 (1) (4@?

by the condition of strict diagonal dominance (3.13), which then shows that (3.14) holds. That completes the inductive step.

We have thus proved that $|u_j| > c_j$ for all $j = 1, 2, \dots, n$. In particular, we deduce that $u_j \neq 0$ for all $j = 1, 2, \dots, n$; hence the LU factorisation $T = LU$ defined by (3.7) exists. Further,

$$\det(T) = \det(L) \det(U) = \det(U) = u_1 u_2 \dots u_n \neq 0,$$

so T is nonsingular.

The formula (3.7) and the inequalities $|u_j| > c_j$, $j = 1, 2, \dots, n$, now imply that

$$\begin{aligned} |u_j| &= |b_j + l_{j,j-1} c_{j-1}| \\ &= |b_j + a_{j,j-1} / u_{j-1}| \\ &> a_{j,j-1}, \quad j = 1, 2, \dots, n, \end{aligned} \tag{3.15}$$

so the elements u_j cannot grow large, and rounding errors are kept under control without pivoting. □

It is easy to see that the same result holds under the weaker assumption that the matrix is diagonally dominant, but not necessarily strictly diagonally dominant, provided that we also require that all the elements c_j , $j = 1, 2, \dots, n-1$, are nonzero (see Exercise 5).

Note also that the matrix constructed in Example 3.2 satisfies this condition, provided that the function r is nonnegative; this often holds in practical boundary value problems.

If the matrix $T \in \mathbb{R}^{n \times n}$ is symmetric and positive definite, as well as tridiagonal, it can be factorised in the form $T = LL^T$, where $L \in \mathbb{R}^{n \times n}$ is lower triangular with nonzero elements only on and immediately below the diagonal. If we use the notation $d_i = |l_{i,i}|$, $e_i = |l_{i,i-1}|$ we easily find from (3.2) and (3.3) that the elements can be calculated in succession from the following formulae:

$$\begin{aligned} d_1 &= b_1, \\ e_i &= c_{i-1} / d_{i-1}, \quad d_i = b_i - e_i^2, \quad i = 2, 3, \dots, n. \end{aligned}$$

This calculation involves about $4n$ operations. Including also the work required by the forward and backsubstitution stages, the complete solution of $Tx = b$ will be found to involve about $10n$ operations. For the tridiagonal matrix the Cholesky factorisation method thus requires more work for the complete solution than the Thomas algorithm; in this case there is no particular advantage in exploiting the symmetry of the matrix in this way.

$$2 \begin{matrix} * & 5 \\ + & \\ ! & \end{matrix} > \quad \&\% \quad :$$

More generally, a system of equations may often involve a matrix of band type.

Definition 3.4 $B^{n \times n}$ is a **band matrix** if there exist nonnegative integers $p < n$ and $q < n$ such that $b_{ij} = 0$ for all $i, j = 1, 2, \dots, n$ such that $p < i - j$ or $q < j - i$. The band is of width $p + q + 1$, with p elements to the left of the diagonal and q elements to the right of the diagonal, in each row. Such a matrix is said to be $\text{Band}(p, q)$.

Thus, for example, a tridiagonal matrix is $\text{Band}(1,1)$, and an $n \times n$ lower triangular matrix is $\text{Band}(n-1,0)$.

An example of a $\text{Band}(1,2)$ matrix $A^{n \times n}$ is shown in Figure 3.1, where each nonzero element in the matrix is identified by an asterisk. In addition to its main diagonal, the matrix has nonzero elements on its lower subdiagonal and two of its superdiagonals.

It is easy to see that, provided that no interchanges are necessary, such a band matrix can be written in the form $B = LU$, where L is $\text{Band}(p,0)$ and U is $\text{Band}(0,q)$ (see Exercise 7). It is also fairly simple to count the operations required in this calculation; the result is approximately proportional to $n^2(p + 2q)$ when n is moderately large. The most common situation has $q = p$, and then the number of operations is approximately proportional to n^2 . As in the tridiagonal case, this is much smaller than n^3 when p and q are fairly small compared with n .

3.4 Monotone matrices

If a positive real number a is increased by $\Delta > 0$ to $a + \Delta$, then its reciprocal a^{-1} decreases to $(a + \Delta)^{-1}$. It is not usually true, however,

that if we increase some or all of the elements of a nonsingular matrix A^{-1} , then the elements of the inverse A will decrease. This useful property holds for the class of monotone matrices defined below.

The discussion in this section is not related to Gaussian elimination and LU factorisation, but it is of relevance in the iterative solution of systems of linear equations with monotone matrices which arise in the course of numerical approximation of boundary value problems for certain ordinary and partial differential equations.

Definition 3.5 *The nonsingular matrix A^{-1} is said to be **monotone** if all the elements of the inverse A are nonnegative.*

Example 3.3 *Suppose that a and d are positive real numbers, and b and c are nonnegative real numbers such that $ad > bc$. Then,*

$$A = \begin{pmatrix} a & b \\ c & d \end{pmatrix}$$

is a monotone matrix. This is easily seen by considering the inverse of the matrix A ,

$$A^{-1} = \frac{1}{ad - bc} \begin{pmatrix} d & b \\ c & a \end{pmatrix},$$

and noting that all elements of A^{-1} are nonnegative.

Next we introduce the concept of ordering in \mathbb{R}^n and $\mathbb{R}^{n \times n}$.

Definition 3.6 *For vectors x and y in \mathbb{R}^n we use the notation*

$$x \leq y, \quad i = 1, 2, \dots, n.$$

In the same way, for matrices A and B in $\mathbb{R}^{n \times n}$ we write

$$A \leq B$$

to mean that

$$a_{ij} \leq b_{ij}, \quad i, j = 1, 2, \dots, n.$$

The sign \leq is read ‘succeeds or is equal to’ or, simply, ‘is greater than or equal to’.

Note that, given two arbitrary matrices A and B in $\mathbb{R}^{n \times n}$, in general none of $A \leq B$, $A = B$ and $B \leq A$ will be true. Therefore the relation $A \leq B$ is a partial, rather than a total, ordering on $\mathbb{R}^{n \times n}$; the same is true of the ordering $A \geq B$.

Theorem 3.5 (i) Suppose that the nonsingular matrix $A \in \mathbb{R}^{n \times n}$ is monotone, $A \geq 0$, and the vectors x and y in \mathbb{R}^n are the solutions of

$$Ax = x, \quad Ay = y,$$

respectively. If $A \geq 0$, then $x \geq y$.

(ii) Suppose that A and B are nonsingular matrices in $\mathbb{R}^{n \times n}$ and that both are monotone. If $A \leq B$, then $B^{-1} \leq A^{-1}$.

Proof (i) Since the elements of A^{-1} are nonnegative and

$$x = A^{-1} (Ax),$$

the result follows from the fact that all elements of the vector $A^{-1} (Ax)$ appearing on the right-hand side of this equality are nonnegative.

(ii) Since $A \leq B$ and all the elements of B^{-1} are nonnegative, it follows that

$$B^{-1}A \leq B^{-1}B = I.$$

In the same way, since all the elements of A^{-1} are nonnegative, it follows that

$$B^{-1} = B^{-1}AA^{-1} \leq A^{-1},$$

as required. □

The following theorem will be useful in Chapter 13.

Theorem 3.6 Suppose that $n \geq 3$ and $T \in \mathbb{R}^{n \times n}$ is a tridiagonal matrix of the form (3.4) with the properties

$$a_i < 0, \quad i = 2, 3, \dots, n, \quad c_i < 0, \quad i = 1, 2, \dots, n-1,$$

and

$$a_i + b_i + c_i = 0, \quad i = 1, 2, \dots, n,$$

where we have followed the convention that $a_1 = 0$, $c_n = 0$; then, the matrix T is monotone.

Proof Let $k = 1, 2, \dots, n$. Column k of the inverse T^{-1} is the solution of the linear system $Tx = e_k$, where e_k is column k of the identity matrix of size n , having a single nonzero element, 1, in row k . By applying the Thomas algorithm to this linear system, it is easy to deduce by induction from (3.7) that $l_j = 0$, $u_j = 0$ and $v_j = 0$ for all j ; the argument is very similar to the proof of Theorem 3.4. It then follows from (3.8) and (3.9) that, in the notation of the Thomas algorithm, the vectors l and u have nonnegative elements. Hence column k of the inverse T^{-1} has nonnegative elements. Since the same is true for each $k = 1, 2, \dots, n$, it follows that T is monotone. \square

3.5 Notes

Symmetric systems of linear algebraic equations arise in the numerical solution of self-adjoint boundary value problems for differential equations with real-valued coefficients.

For further details on the Cholesky factorisation, the reader may consult any of the books listed in the Notes at the end of Chapter 2, particularly Chapter 10 of N.J. Higham, *Accuracy and Stability of Numerical Algorithms*, SIAM, Philadelphia, 1996.

Classical iterative methods for the solution of systems of linear equations with monotone matrices are discussed, for example, in

% & , *Matrix Iterative Analysis*, Prentice-Hall, Englewood Cliffs, NJ, 1962.

A more recent reference on iterative algorithms for linear systems is

' () * " *Iterative Solution Methods*, Cambridge University Press, Cambridge, 1996.

In particular, Chapter 6 of Axelson's book considers the relevance of monotone matrices in the context of iterative solution of systems of linear equations.

Theorem 3.6 is a slight variation on the following general result.

Theorem 3.7 *A sufficient condition for A^{-1} to be a monotone matrix is that A is an M -matrix, that is, (a) $a_{ii} > 0$ for all $i, j = 1, 2, \dots, n$ such that $i = j$, and (b) there exists a vector w with positive elements such that all elements of $A^{-1}w$ are positive.*

Exercises

3.1 Find the Cholesky factorisation of the matrix

$$A = \begin{pmatrix} 4 & 6 & 2 \\ 6 & 10 & 3 \\ 2 & 3 & 5 \end{pmatrix}.$$

3.2 Use the method of Cholesky factorisation to solve the system of equations

$$\begin{aligned} x^2 + 2x + 2x &= 4, \\ 2x^2 + 5x + 3x &= 7, \\ 2x^2 + 3x + 6x &= 10. \end{aligned}$$

3.3 Let $n \geq 3$. The $n \times n$ tridiagonal matrix T has the diagonal elements

$$T_{ii} = 2, \quad i = 1, 2, \dots, n,$$

and the off-diagonal elements

$$T_{i,i-1} = T_{i-1,i} = 1, \quad i = 1, 2, \dots, n-1.$$

In the factorisation $T = LU$, where L is unit lower triangular and U is upper triangular, show that

$$L_{ii} = i/(i+1), \quad i = 1, 2, \dots, n,$$

and find expressions for the elements of U . What is the determinant of T ?

3.4 Let $n \geq 3$ and $1 \leq k \leq n$. Define the vector v with elements given by

$$v_i = \begin{cases} i(n+1-k), & i = 1, \dots, k, \\ k(n+1-i), & i = k+1, \dots, n. \end{cases}$$

Evaluate M_{ij} , the inner product of the vector v with column j of the matrix T defined in Exercise 3. (The inner product (v, w) of two vectors v and w in \mathbb{R}^n is defined as the real number $\sum_{i=1}^n v_i w_i$.) Hence give expressions for the elements of the inverse matrix T^{-1} , and verify that this inverse is symmetric. Find the ∞ -norm of the inverse, $\|T^{-1}\|_\infty$, and show that the condition number of T is

$$\kappa(T) = \frac{1}{2}(n+1), \quad n \text{ odd}.$$

What is the condition number $\kappa(T)$ when n is even?

3.5 Given that $n \geq 3$, in the notation of Theorem 3.4 suppose that

$$b_j = a_j + c_j, \quad j = 1, 2, \dots, n,$$

and

$$c_j > 0, \quad j = 1, 2, \dots, n-1,$$

with the convention that $a_n = 0$ and $c_n = 0$. Show that the factorisation $T = LU$ exists without pivoting, and can be constructed by the Thomas algorithm. Give an example of a matrix T which satisfies these conditions, except that $c_k = 0$ for some $k = 1, 2, \dots, n-1$ and such that T is singular and cannot be written in the form $T = LU$ without pivoting.

3.6 Let $n \geq 3$ and suppose that the matrix $T \in \mathbb{R}^{n \times n}$ is tridiagonal. Show that there exists a permutation matrix $P \in \mathbb{R}^{n \times n}$ such that

$$PA = LU$$

where $L \in \mathbb{R}^{n \times n}$ is unit lower triangular with at most two nonzero elements in each row, and $U \in \mathbb{R}^{n \times n}$ is upper triangular with at most three nonzero elements in each row.

3.7 Suppose that the matrix B is $\text{Band}(p, q)$, and that there exists a factorisation $B = LU$ without row interchanges. Show that L is $\text{Band}(p, 0)$ and U is $\text{Band}(0, q)$.

3.8 Suppose that $n \geq 4$, that the matrix $A \in \mathbb{R}^{n \times n}$ is $\text{Band}(3, 3)$, and has the LU factorisation $A = LU$, so that $L \in \mathbb{R}^{n \times n}$ is $\text{Band}(3, 0)$ and $U \in \mathbb{R}^{n \times n}$ is $\text{Band}(0, 3)$. Suppose also that $a_{i,i} = 0, a_{i,i+1} = 0$ for $i = 1, 2, \dots, n-2$. By considering u_{ij} and l_{ij} , or otherwise, show that in general the elements $l_{i,i-1}$ and $u_{i,i+2}$ are not zero.

Simultaneous nonlinear equations

4.1 Introduction

In Chapter 1 we discussed iterative methods for the solution of a single nonlinear equation of the form $f(x) = 0$ where f is a continuous real-valued function of a single real variable. In Chapters 2 and 3, on the other hand, we were concerned with direct (as opposed to iterative) methods for systems of linear equations. The purpose of the present chapter is to extend the techniques developed in Chapter 1 to systems of simultaneous nonlinear equations for functions of several real variables. We shall concentrate on two methods: the generalisation of simple iteration, usually referred to as simultaneous iteration, and Newton's method.

Given that $x = (x_1, \dots, x_n)$, as in Chapters 2 and 3 we denote by $\|x\|$ the ∞ -norm of x defined by

$$\|x\| = \max |x_i|.$$

Throughout the chapter, \mathbb{R}^n will be thought of as a linear space equipped with the ∞ -norm; with only minor alterations all of our results can be restated in the p -norm with $p \in [1, \infty)$ on replacing $\|x\|$ by $\|x\|_p$ throughout. We begin with some basic definitions which involve the concept of *open ball* defined in Section 2.7.

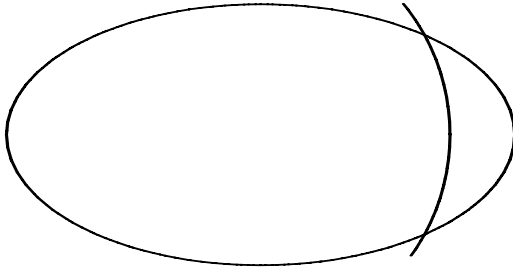
Let $x_0 \in \mathbb{R}^n$; the open ball in \mathbb{R}^n (with respect to the ∞ -norm) of radius $r > 0$ and centre x_0 is defined as the set

$$B(x_0, r) = \{x \in \mathbb{R}^n : \|x - x_0\| < r\}.$$

A set $D \subseteq \mathbb{R}^n$ is said to be an **open set** in \mathbb{R}^n if for every $x \in D$ there exists $\delta = \delta(x) > 0$ such that $B(x, \delta) \subseteq D$ (see Figure 4.1). For example, any open ball in \mathbb{R}^n is an open set in \mathbb{R}^n . Given $x_0 \in \mathbb{R}^n$, any open set



$\epsilon > 0$ such that $B(x, \epsilon) \cap D = \emptyset$. As $(x_n) \subset D$, no member of the sequence (x_n) can enter $B(x, \epsilon)$. This, however, contradicts the fact that (x_n) converges to x . The contradiction implies that D is closed. \square



Example 4.2 Let us suppose that $A \in \mathbb{R}^{n \times n}$ and $f: D \rightarrow \mathbb{R}^n$. On letting $(x) = A^{-1} f(x)$ we deduce that the problem of solving the system of simultaneous linear equations considered in Chapters 2 and 3 can be restated in the form: find (x) such that $(x) = 0$.

Let us assume that we have transformed the equation $(x) = 0$ into an equivalent form $(x) = g(x)$, where $g: D \rightarrow D$ is a continuous function, defined on the closed subset D of \mathbb{R}^n , such that $g(D) \subset D$. For example, one can choose $(x) = A^{-1} f(x)$, with α a suitable parameter. By ‘equivalent’ we mean that $(x) \in D$ satisfies $(x) = 0$ if, and only if, $(x) = g(x)$. Any $(x) \in D$ such that $(x) = g(x)$ is called a **fixed point** of the function g in D . Thus the problem of finding a solution $(x) \in D$ to the equation $(x) = 0$ has been converted into one of finding a fixed point in D of the function g . We embark on the latter task by considering the natural extension to \mathbb{R}^n of the simple iteration discussed in Section 1.2 for the solution of the scalar nonlinear equation $g(x) = x$.

Definition 4.1 Suppose that $g: D \rightarrow D$ is a function, defined and continuous on a closed subset D of \mathbb{R}^n , such that $g(D) \subset D$. Given that $(x_0) \in D$, the recursion defined by

$(x_k) = g(x_{k-1}), \quad k = 0, 1, 2, \dots, \quad (4.3)$

is called a **simultaneous iteration**. For $n = 1$ the recursion (4.3) is just the simple iteration considered in (1.3).

Note that here we use the superscript k as the sequence index; following the convention adopted in Chapters 2 and 3, we reserve subscripts for labelling the entries of vectors. Thus $x_i^{(k)}$ is entry i of the vector (x_k) , the k th member of the sequence (x_k) . The motivation behind the definition of the simultaneous iteration (4.3) is, of course, our hope that, under suitable conditions on g and D , the sequence (x_k) will converge to a fixed point (x^*) of g .

Two remarks are in order at this point. First, it is easy to show that if a sequence of vectors (x_k) converges in \mathbb{R}^n to (x) in the norm $\| \cdot \|_p$, then it also converges to this same limit in the norm $\| \cdot \|_q$ for any $p, q \in [1, \infty)$. To see this, note that

$$\|x_k - x\|_q \leq \|x_k - x\|_p, \quad (4.4)$$

for $1 < p < \infty$, and take $\epsilon = \frac{1}{k}$ to deduce that, as $k \rightarrow \infty$, convergence in the ∞ -norm implies convergence in the p -norm for any $p \in [1, \infty)$, and *vice versa*.

Any function f that satisfies a Lipschitz condition on a set D is continuous on D . For let $x, y \in D$ and $\epsilon > 0$; then, on defining $\delta = \epsilon/L$, we deduce from (4.5) that if $\|x - y\| < \delta$ for some $x, y \in D$, then

$$\|f(x) - f(y)\| \leq L \|x - y\| < \epsilon.$$

It follows from (4.4) that if f satisfies a Lipschitz condition on D in the ∞ -norm then it also does so in the p -norm for any $p \in [1, \infty)$, and *vice versa*. However, in general, the size of the constant L may depend on the choice of norm. Specifically, if f is a contraction on a set D in the ∞ -norm (*i.e.*, (4.5) holds with $L < 1$), then f need not be a contraction in the p -norm, unless $L < n^{-1/p}$. (See Exercise 1.) Conversely, if f is a contraction on D in the p -norm for some $p \in [1, \infty)$, it does not follow that f is a contraction on D in the ∞ -norm.

For example, suppose that $f : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ is the linear function defined by $f(x) = Ax$, where A is the 2×2 matrix

$$A = \begin{pmatrix} 3/4 & 1/3 \\ 0 & 3/4 \end{pmatrix}.$$

This function f satisfies a Lipschitz condition on \mathbb{R}^2 in ∞ -norm for any $p \in [1, \infty)$, and if L is a Lipschitz constant for f in the p -norm, then $L \leq \|A\|_p$, in the subordinate matrix norm. It is easy to see that $\|A\|_1 = \|A\|_\infty = 13/12$, and a small calculation gives $\|A\|_p = 0.935$ to three decimal digits. Hence the function f is a contraction in the 2-norm, but not in the 1- or ∞ -norm.

Our next result is a direct generalisation of Theorem 1.3 formulated in Chapter 1.

Theorem 4.1 (Contraction Mapping Theorem) *Suppose that D is a closed subset of \mathbb{R}^n , $f : D \rightarrow \mathbb{R}^n$ is defined on D , and $f(D) \subset D$. Suppose further that f is a contraction on D in the ∞ -norm. Then, f has a unique fixed point in D , and the sequence (x_n) defined by (4.3) converges to x^* for any starting value $x_0 \in D$.*

Proof Assuming that f has a fixed point x^* in D , the uniqueness of the fixed point is easy to show: for suppose that y^* is also a fixed point of f in D . Then, by (4.5),

$$\|x^* - y^*\| = \|f(x^*) - f(y^*)\| \leq L \|x^* - y^*\|,$$

i.e., $(1 - L) \|\mathbf{x} - \mathbf{0}\| \leq L \|\mathbf{x} - \mathbf{0}\|$. Since $L < 1$, and $\|\cdot\|$ is a norm, it follows that $\mathbf{x} = \mathbf{0}$, and hence $\mathbf{x} = \mathbf{0}$. Consequently, if \mathcal{D} has a fixed point in \mathcal{D} , then this is the unique fixed point of T in \mathcal{D} .

Now, still *assuming* that T possesses a fixed point \mathbf{x}^* in \mathcal{D} , we shall show that the sequence (\mathbf{x}_k) defined by (4.3) converges to \mathbf{x}^* for any starting value $\mathbf{x}_0 \in \mathcal{D}$.

for all $k \geq 2$. We then deduce by induction that

$$\|x^{(k)} - x^{(k-1)}\| \leq L^{k-1} \|x^{(1)} - x^{(0)}\|, \quad k \geq 1. \quad (4.8)$$

Suppose that m and k are positive integers and $m \geq k + 1$. Then, by repeated application of the triangle inequality in the ∞ -norm and using (4.8), we have that

$$\begin{aligned} \|x^{(m)} - x^{(k)}\| &= (\|x^{(m)} - x^{(m-1)}\| + \|x^{(m-1)} - x^{(m-2)}\| + \dots + \|x^{(k+1)} - x^{(k)}\|) \\ &\leq (L^{m-1} + L^{m-2} + \dots + L^k) \|x^{(1)} - x^{(0)}\| \\ &= L^k (L^{m-k} + L^{m-k-1} + \dots + 1) \|x^{(1)} - x^{(0)}\| \\ &= L^k \frac{1 - L^{m-k+1}}{1 - L} \|x^{(1)} - x^{(0)}\|, \end{aligned} \quad (4.9)$$

where, in the transition to the last line, we made use of the fact that the geometric series $1 + L + L^2 + \dots$, with $L \in (0, 1)$, sums to $1/(1 - L)$.

As $\lim_{m \rightarrow \infty} L^{m-k+1} = 0$, it follows from (4.9) that $(x^{(k)})$ is a Cauchy sequence in \mathbb{R}^n ; that is, for each $\epsilon > 0$ there exists $k = k(\epsilon)$ (defined by (4.6) above) such that

$$\|x^{(m)} - x^{(k)}\| < \epsilon, \quad m, k \geq k(\epsilon). \quad (4.10)$$

Any Cauchy sequence in \mathbb{R}^n is convergent in \mathbb{R}^n ; consequently, there exists x^* such that $x^* = \lim_{m \rightarrow \infty} x^{(m)}$. Further, since F satisfies a Lipschitz condition on D , the discussion in the paragraph following Definition 4.2 shows that F is continuous on D . Hence, by Lemma 4.2,

$$F(x^*) = \lim_{m \rightarrow \infty} F(x^{(m)}) = \lim_{m \rightarrow \infty} 0 = 0,$$

which proves that x^* is a fixed point of T .

It remains to show that $x^* \in D$. This follows from Lemma 4.1 since $(x^{(k)}) \subset D$, $x^* = \lim_{m \rightarrow \infty} x^{(m)}$ and D is closed. \square

As a byproduct of the proof, we deduce from (4.7) that, given a positive tolerance ϵ , one can compute an approximation $x^{(k)}$ to the unknown solution x^* using (4.3) in no more than $k = k(\epsilon)$ iterations so that the approximation error $\|x^{(k)} - x^*\|$, measured in the ∞ -norm, is less than ϵ ; the integer $k(\epsilon)$ is defined by (4.6).

The next theorem relates the constant L from the Lipschitz condition (4.5) to the partial derivatives of F , giving a more practically useful sufficient condition for convergence.

Definition 4.3 Let $\mathbf{g} = (g_1, \dots, g_n) : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be a function defined and continuous in an (open) neighbourhood $N(\mathbf{a})$ of \mathbf{a} . Suppose further that the first partial derivatives $\frac{\partial g_i}{\partial x_j}$, $j = 1, \dots, n$, of g_i exist at \mathbf{a} for $i = 1, \dots, n$. The **Jacobian matrix** $J(\mathbf{g})(\mathbf{a})$ of \mathbf{g} at \mathbf{a} is the $n \times n$ matrix with elements

$$J(\mathbf{g})(\mathbf{a}) = \left(\frac{\partial g_i}{\partial x_j}(\mathbf{a}) \right)$$

for $i = 1, \dots, n$. Now $x_j = y_j$ for all $j = 1, \dots, n$, and so (4.12) gives

$$g_i(x) = g_i(y) = \frac{g_i}{x_i} (x_i + (1 - x_i) y_i),$$

$$J_i(x + (1 - t)y) =$$

for all $i = 1, \dots, n$. Consequently, for any $t \in [0, 1]$, $\bar{B}(x, t)$,

$$\|g(x + (1 - t)y) - g(y)\| \leq \max_i J_i(t + (1 - t)x) \|x - y\| \leq -(1 + K) \|x - y\|, \tag{4.13}$$

due to (4.11), given that $t + (1 - t)x \in \bar{B}(x)$ for all $t \in [0, 1]$. It follows that g satisfies a Lipschitz condition (4.5), in the ∞ -norm, on the closed ball $\bar{B}(x)$ with $L = -(1 + K) < 1$. Furthermore, on selecting $t = 1$ in (4.13) we get that

$$\|g(x) - g(y)\| \leq \|g(x) - g(y)\| <$$

for all $x, y \in \bar{B}(x)$. Hence, $(\bar{B}(x)) \subset \bar{B}(x)$. The convergence of the iteration (4.3) to x , for an arbitrary starting value $x_0 \in \bar{B}(x)$, now follows from Theorem 4.1. \square

We close this section with an example which illustrates the application of the method of simultaneous iteration to the solution of a system of nonlinear equations.

Example 4.4 *Let us consider, as in Example 4.1, the system of two simultaneous nonlinear equations in the unknowns x and y , defined by*

$$\begin{aligned} x + x^2 - 1 &= 0, \\ 5x + 21x^2 - 9 &= 0. \end{aligned}$$

Here $x = (x, x)$ and $f = (f, f)$ with

$$\begin{aligned} f_1(x, x) &= x + x^2 - 1, \\ f_2(x, x) &= 5x + 21x^2 - 9. \end{aligned}$$

Let us suppose that we need to find the solution of the system $f(x) = 0$ in the first quadrant of the (x, x) -coordinate system.

Of course, the example is a little artificial, since we already know from Example 4.1 that $x = (\sqrt{3}/2, 1/2)$ is the required solution. In what follows, however, we proceed as if we knew nothing about the location

of \mathbb{R}^2 . Our aim here is to illustrate the construction of the function from \mathbb{R}^2 and the verification of the hypotheses of Theorem 4.1.

Let us rewrite the two equations as

$$x = 1 - x^2, \quad x = \frac{1}{21} (9 - 5x^2),$$

and define $g(x, x)$ and $g(x, x)$ as the right-hand sides of these, respectively. We consider the simultaneous iteration

for all x and y in D . Therefore, also,

$$\|J(x, y)\| \leq L$$

with

$$L = \max \|J(x, y)\| < 1. \tag{4.15}$$

With our choice of D , (4.15) holds with $L = \max \|J(x, y)\| = 0.75 < 1$. Furthermore, it is easy to check that $D \subset D$. Thus we deduce from Theorem 4.1 that T has a unique fixed point in D – we call this fixed point (x^*, y^*) , for the sake of consistency with the notation in Example 4.1; moreover, the sequence (x_k, y_k) defined by (4.14) converges to (x^*, y^*) .

After all these preparations you are now probably curious to see what the successive iterates look like: Table 4.1 gives a flavour of the behaviour of the sequence (x_k, y_k) , with the starting value chosen as $(x_0, y_0) = (0.5, 0.3)$. You can see that after 15 iterations the first 5 decimal digits have settled to their correct values.

4.3 Relaxation and Newton’s method

We now go on to apply the ideas developed in the previous section to the construction of an iteration which converges to a solution of the equation $F(x, y) = 0$, where $F: D \rightarrow \mathbb{R}^2$. One way of constructing such a sequence is by relaxation.

Definition 4.4 *The recursion*

$$x_{k+1} = \alpha x_k + (1 - \alpha) G(x_k, y_k), \quad y_{k+1} = \beta y_k + (1 - \beta) H(x_k, y_k), \quad k = 0, 1, 2, \dots, \tag{4.16}$$

where G, H is given and where $\alpha, \beta \in [0, 1]$ is a constant, is called **simultaneous relaxation**.

Suppose that the sequence (x_k, y_k) converges to a limit (x^*, y^*) and F is continuous in a neighbourhood of (x^*, y^*) ; then, on passing to the limit $k \rightarrow \infty$ in (4.16), we deduce that (x^*, y^*) is a solution of the equation $F(x, y) = 0$.

Simultaneous relaxation is evidently a simultaneous iteration defined by taking $(x_{k+1}, y_{k+1}) = G(x_k, y_k)$.

1	!	6	,	1	!	G	2	+	4
G	6	B	4	7	D	1	B	8	D
1	"	5	5	0					

Table 4.1. The first 15 iterates in the sequence $x_k = (x_k^1, x_k^2)$ defined by (4.14), with starting value $(0.5, 0.3)$. The exact solution is $x^* = (\sqrt{3}/2, 1/2) = (0.866025403784439, 0.500000000000000)$ to 15 decimal digits.

k	x_k^1	x_k^2
0	0.5	0.3
1	0.707106781186548	0.707106781186548
2	0.866025403784439	0.500000000000000
3	0.866025403784439	0.500000000000000
4	0.866025403784439	0.500000000000000
5	0.866025403784439	0.500000000000000
6	0.866025403784439	0.500000000000000
7	0.866025403784439	0.500000000000000
8	0.866025403784439	0.500000000000000
9	0.866025403784439	0.500000000000000
10	0.866025403784439	0.500000000000000
11	0.866025403784439	0.500000000000000
12	0.866025403784439	0.500000000000000
13	0.866025403784439	0.500000000000000
14	0.866025403784439	0.500000000000000
15	0.866025403784439	0.500000000000000

Theorem 4.3 Suppose that $f(x) = 0$, and that all the first partial derivatives of $f = (f_1, \dots, f_n)$ are defined and continuous in some (open) neighbourhood of x^* , and satisfy a condition of strict diagonal dominance at x^* ; i.e.,

$$\frac{\partial f_i}{\partial x_i}(x^*) > \sum_{j \neq i} \left| \frac{\partial f_i}{\partial x_j}(x^*) \right|, \quad i = 1, 2, \dots, n. \tag{4.17}$$

Then, there exist $\delta > 0$ and a positive constant α such that the relaxation iteration (4.16) converges to x^* for any x_0 in the closed ball $\bar{B}(x^*, \delta)$ of radius δ , centre x^* .

Proof The elements of the Jacobian matrix $J(x) = \left(\frac{\partial f_i}{\partial x_j} \right)_{i,j=1}^n$ of the function $f(x) = (f_1(x), \dots, f_n(x))$ at $x = x^*$ are

$$J_{ij}(x^*) = \begin{cases} 1 - \frac{\partial f_i}{\partial x_i}(x^*), & j = i, \\ -\frac{\partial f_i}{\partial x_j}(x^*), & j \neq i, \end{cases} \quad i, j = 1, \dots, n.$$

We now define

$$m = \max \frac{f}{x}(\cdot)$$

and then choose $\epsilon = 1/m$. Under hypothesis (4.17), $m > 0$ and therefore $\epsilon > 0$. This choice of ϵ ensures that all the diagonal elements $J_{ii}(\cdot)$, $i = 1, \dots, n$, of $J(\cdot)$ are nonnegative. Moreover, for any $i = 1, \dots, n$,

the system of equations $() = \mathbf{0}$. It is implicitly assumed that the matrix $J()$ exists and is nonsingular for each $k = 0, 1, 2, \dots$.

The next theorem is concerned with the convergence of Newton's method. As in the scalar case, for a starting value that is sufficiently close to the solution of $() = \mathbf{0}$, Newton's method converges quadratically. The precise definition of quadratic convergence is given below: it resembles Definition 1.7 of Chapter 1.

Definition 4.6 Suppose that $()$ is a convergent sequence in and $\lim =$. We say that $()$ converges to **with at least order $q > 1$** , if there exist a sequence $()$ of positive real numbers converging to 0, and $\mu > 0$, such that

$$, \quad k = 0, 1, 2, \dots, \quad \text{and} \quad \lim \text{---} = \mu. \quad (4.19)$$

If (4.19) holds with $=$, $k = 0, 1, 2, \dots$, then the sequence $()$ is said to converge to **with order q** . In particular, if $q = 2$, then we say that the sequence $()$ converges to **quadratically**.

Again, due to (4.4), if a sequence $()$ converges quadratically in the $-$ -norm, then it also does so in the p -norm for any $p \in [1,)$, though the constant μ may be different.

Theorem 4.4 Suppose that $() = \mathbf{0}$, that in some (open) neighbourhood $N()$ of $,$ where $is defined and continuous, all the second-order partial derivatives of are defined and continuous, and that the Jacobian matrix $J()$ of at the point is nonsingular. Then, the sequence $()$ defined by Newton's method (4.18) converges to the solution provided that is sufficiently close to ; the convergence of the sequence $()$ to is at least quadratic.$

Proof Let us begin by writing Newton's method as a simultaneous iteration $= (), k = 0, 1, 2, \dots$, as in (4.3), with given and

$$() = [J()]^{-1} ().$$

The idea of the proof is to verify that the function satisfies all the conditions of Theorem 4.2 in a certain closed ball centred at $,$ the fixed point of $,$ and thus deduce that the sequence $()$ converges to $.$

As the function $\det J()$ is continuous in $N()$ and $\det J() = 0$, there exists > 0 such that $\det J() \neq 0$ for all $\bar{B}() \subset N().$

Further, as the entries of $[J(\cdot)]^{-1}$ depend continuously on the entries of $J(\cdot)$ and since the entries of $J(\cdot)$ are continuous functions of \cdot in $N(\cdot)$, we deduce that $[J(\cdot)]^{-1}(\cdot)$ is a continuous function on $\bar{B}(\cdot)$; therefore,

$$(\cdot) = [J(\cdot)]^{-1}(\cdot)$$

is also a continuous function on $\bar{B}(\cdot)$. For later reference, we note that $[J(\cdot)]^{-1}$

and so

$$-n \mathbf{A} \mathbf{C} \quad .$$

On writing $\mathbf{M} = -n \mathbf{A} \mathbf{C}$, we then deduce by induction that

$$\frac{1}{M} \mathbf{M} \quad , \quad k = 0, 1, 2, \dots .$$

Suppose that $\bar{\mathbf{B}}(\cdot)$ where $-\min(1, 1/M)$. Then,

$$\mathbf{M} \quad \frac{1}{2}, \quad k = 0, 1, 2, \dots ,$$

and hence

$$\frac{1}{M} \quad \frac{1}{2}$$

This implies that convergence is at least quadratic (on choosing $= M^{-2}$ and $q = 2$ in Definition 4.6). □

Newton's method is defined in (4.18) by using the inverse of the Jacobian matrix. As we saw in Chapter 2 it is more efficient to avoid inverting a matrix, if possible. In practice the method is therefore implemented by writing (4.18) in the form

$$\mathbf{J}(\cdot) [\quad] = (\cdot) . \tag{4.23}$$

Given the vector \cdot , we calculate (\cdot) and the Jacobian matrix $\mathbf{J}(\cdot) \times$, and then solve the system of linear equations (4.23) by Gaussian elimination; this gives the increment vector \cdot , which is added to \cdot to obtain the new iterate \cdot .

Example 4.5 *We close this section with an example which illustrates the application of Newton's method. Consider the simultaneous nonlinear equations*

$$\begin{aligned} f(x, y, z) &= x + y + z - 1 = 0, \\ f(x, y, z) &= 2x + y - 4z = 0, \\ f(x, y, z) &= 3x - 4y + z = 0. \end{aligned}$$

Letting $\mathbf{f} = (f, f, f)$ and $\mathbf{x} = (x, y, z)$, the aim of the exercise is to determine the solution to the equation $\mathbf{f}(\mathbf{x}) = \mathbf{0}$ contained in the first octant $(x, y, z) : x > 0, y > 0, z > 0$ in \cdot .



P are, very roughly, $y = 0.5$ and $z = 0.5$, it is reasonable to choose as starting value for the Newton iteration the point

Chapter 1, and also an infinite number of complex solutions. It is easy to see from the periodic character of $e^{\#}$ that the equation has a solution near $w = (2m + \pi)i$, $v = \overline{-1}$, for integer values of m ; a better estimate is given in Exercise 9. It is a good deal more difficult to prove that there are no other solutions.

The behaviour of Newton's method for this problem may be illustrated by showing a picture of the complex plane, with the sets S depicted in different colours. In our example we cannot, of course, show more than a small number of the solutions, and cannot use an infinite number of colours. We have therefore coloured the sets with six colours cyclically, so that, for example, the sets S_1, S_2, S_3, \dots have the same colour. The background colour, white, represents the set S of points from which the iteration converges to the real negative root. It includes most of the negative half-plane. Successive pictures in the series from Figure 4.5 to Figure 4.9 show a magnified view of a small region of the previous picture, the region being outlined in black. In Figure 4.4 the black crosses mark the positions of solutions of $f(z) = 0$. The pictures show in a striking way the fractal behaviour of the boundary of a set. Figure 4.9 is very similar to Figure 4.5; the former is a magnified view of a small part of Figure 4.5, with a magnification of about 50000 in each direction. The same sort of behaviour is repeated when the picture is magnified indefinitely.

4.5 Notes

For an introduction to the topology of \mathbb{C} , including the definitions of open set, closed set, continuity, convergence and Cauchy sequence, the reader is referred to any standard textbook on the subject; see, *e.g.*,

% *Principles of Mathematical Analysis*, Third Edition, International Series in Pure and Applied Mathematics, McGraw-Hill, New York, Auckland, Düsseldorf, 1976,

&) *Introduction to Mathematical Analysis*, Addison-Wesley, Reading, MA, 1996.

Our first remark concerns the Contraction Mapping Theorem, Theorem 4.1, which is a direct generalisation of Theorem 1.3 from Chapter 1. Comparing the proofs of Theorems 1.3 and 4.1, we see that the proof of Theorem 1.3 is much simpler. This is not accidental: in the case of a single equation $x = g(x)$, involving a real-valued function g of a single real variable x , the existence of a fixed point follows directly from

Theorem 1.2, Brouwer's Fixed Point Theorem on a bounded closed interval of the real line. On the other hand, for the simultaneous system of equations $\mathbf{y} = \mathbf{f}(\mathbf{x})$ in \mathbb{R}^n considered in Theorem 4.1 we had to invoke the completeness of \mathbb{R}^n (i.e., the property that every Cauchy sequence in \mathbb{R}^n is a convergent sequence) to show the existence of a fixed point. An alternative, shorter proof of Theorem 4.1 could have been devised by applying Brouwer's Fixed Point Theorem in \mathbb{R}^n .

Theorem 4.5 (Brouwer's Fixed Point Theorem) *Let us assume that D is a nonempty, closed, bounded and convex subset of \mathbb{R}^n . Suppose further that $f : D \rightarrow \mathbb{R}^n$ is a continuous function defined on D such that $f(D) \subset D$. Then, there exists $\mathbf{x} \in D$ such that $f(\mathbf{x}) = \mathbf{x}$.*

A set $D \subset \mathbb{R}^n$ is said to be convex if, whenever \mathbf{x} and \mathbf{y} belong to D , also

$$\mathbf{z} = t\mathbf{x} + (1-t)\mathbf{y} \in D \quad \text{for } t \in [0, 1].$$

For example, any nonempty interval of the real line $I \subset \mathbb{R}$ is a convex set, as is a nonempty (open or closed) ball in \mathbb{R}^n , $n \geq 2$. Unfortunately, when $n \geq 2$ the proof of Theorem 4.5 is nontrivial and is well beyond the scope of this book.

Benoit Mandelbrot (1924–) has been largely responsible for the present interest in fractal geometry and its connections with iterative methods. Mandelbrot highlighted in his book

! + _____, *Fractals: Form, Chance, and Dimension*, W.H. Freeman, San Francisco, 1977,

and, more fully, in

! + _____ " *The Fractal Geometry of Nature*, W.H. Freeman, New York, 1983,

the omnipresence of fractals both in mathematics and elsewhere in nature. In relation with the subject of this chapter, we note that the **Mandelbrot set** is a connected set of points in the complex plane defined as follows. Choose a point z in the complex plane, and consider the iteration $z_{n+1} = z_n^2 + z_n$, $n = 0, 1, 2, \dots$. If the sequence z_0, z_1, z_2, \dots remains within a distance of 2 from the origin for ever, then the point z

1 9 _____ 1 1 _____ 4> _____ , _____ K
4 ' _____) # _____ 5 + _____ : _____
, #@@&4 _____ " _____

is said to be in the Mandelbrot set. If the sequence diverges from the origin, then the point z is not in the set.

A standard reference for theoretical results concerning the convergence of Newton's method in complete normed linear spaces is

§) - " *Functional Analysis*, Second edition, Pergamon Press, Oxford, New York, 1982.

A further significant book in the area of iterative solution of systems of nonlinear equations is the text by

+ ' % " *Iterative Solution of Non-linear Equations in Several Variables*, Reprint of the 1970 original, Classics in Applied Mathematics, 30, SIAM, Philadelphia, 2000.

It gives a comprehensive treatment of the numerical solution of n nonlinear equations in n unknowns, covering asymptotic convergence results for a number of algorithms, including Newton's method, as well as existence theorems for solutions of nonlinear equations based on the use of topological degree theory and Brouwer's Fixed Point Theorem.

Exercises

4.1 Suppose that the function is a contraction in the ∞ -norm, as in (4.5). Use the fact that

$$\|f(x) - f(y)\|_\infty \leq L \|x - y\|_\infty$$

to show that is a contraction in the p -norm if $L < n^{-1}$.

4.2 Show that the simultaneous equations $(x, x) = \mathbf{0}$, where (f, f) , with

$$f(x, x) = x + x - 25, \quad f(x, x) = x - 7x - 25,$$

have two solutions, one of which is $x = 4, x = 3$, and find the other. Show that the function does not satisfy the conditions of Theorem 4.3 at either of these solutions, but that if the sign of f is changed the conditions are satisfied at one solution, and that if is replaced by $(f - f, f)$, then the conditions are satisfied at the other. In each case, give a value of the relaxation parameter which will lead to convergence.

- 4.3 The complex-valued function $z \mapsto g(z)$ of the complex variable z is holomorphic in a convex region Ω containing the point z_0 , at which $g'(z_0) = \lambda$. By applying the Mean Value Theorem (Theorem A.3) to the function $g(tu + tv)$ of the real variable t defined by $g(tu + tv) = g((1-t)u + tv)$ show that if u and v lie in Ω , then there is a complex number θ in Ω such that

$$g(u) - g(v) = (u - v)g'(\theta).$$

Hence show that if $|g'(z_0)| < 1$, then the complex iteration defined by $z_{k+1} = g(z_k)$, $k = 0, 1, 2, \dots$, converges to z_0 provided that z_0 is sufficiently close to z_0 .

- 4.4 Suppose that in Exercise 3 the real and imaginary parts of g are u and v , so that $g(x + iy) = u(x, y) + iv(x, y)$, $i = \sqrt{-1}$. Show that the iteration defined by $z_{k+1} = g(z_k)$, $k = 0, 1, 2, \dots$, where $z_k = (u(x_k, x_k), v(x_k, x_k))$, generates the real and imaginary parts of the sequence defined in Exercise 3. Compare the condition for convergence given in that exercise with the sufficient condition given by Theorem 4.2.
- 4.5 Verify that the iteration $z_{k+1} = g(z_k)$, $k = 0, 1, 2, \dots$, where $g = (g_1, g_2)$ and g_1 and g_2 are functions of two variables defined by

$$g_1(x, x) = -(x^2 - x + 3), \quad g_2(x, x) = -(2x^2 + 1),$$

has the fixed point $z_0 = (1, 1)$. Show that the function g does not satisfy the conditions of Theorem 4.3. By applying the results of Exercises 3 and 4 to the complex function g defined by

$$g(z) = -(z^2 + 3 + i), \quad z = x + iy, \quad i = \sqrt{-1},$$

show that the iteration, nevertheless, converges.

- 4.6 Suppose that all the second-order partial derivatives of the function $f: \mathbb{R}^n \rightarrow \mathbb{R}$ are defined and continuous in a neighbourhood of the point a in \mathbb{R}^n , at which $\nabla f(a) = \mathbf{0}$. Assume also that the Jacobian matrix, $J_f(a)$, of f is nonsingular at a , and denote its inverse by $K(a)$ at all a for which it exists. Defining the Newton iteration by $a_{k+1} = g(a_k)$, $k = 0, 1, 2, \dots$, with a_0 given, where $g(a) = a - K(a)^{-1} \nabla f(a)$, show that the (i, j) -entry

4.9 Suppose that the equation $e^z = z + 2$, $z \in \mathbb{C}$, has a solution

$$z = (2m + i)t + \ln[(2m + i)^2] + i\pi,$$

where m is a positive integer and $t = \sqrt{m^2 + 1}$. Show that

$$\operatorname{Re} z = \ln[1 - t(\ln(2m + i) + i + 2)/(2m + i)]$$

and deduce that $\operatorname{Re} z \sim (\ln m/m)$ for large m .

(Note that $\ln(1 + it) < t$ for all $t > 0$.)

°

' 5
!

'\$&

'\$%

&

'\$ 5
)(!

'\$&%('

'\$&%%()&

Eigenvalues and eigenvectors of a symmetric matrix

5.1 Introduction

Eigenvalue problems for symmetric matrices arise in all areas of applied science. The terminology *eigenvalue* comes from the German word *Eigenwert* which means proper or characteristic value. The concept of eigenvalue first appeared in an article on systems of linear differential equations by the French mathematician d'Alembert in the course of studying the motion of a string with masses attached to it at various points.

Let us recall from Chapter 2 the definition of eigenvalue and eigenvector.

Definition 5.1 Suppose that A is an $n \times n$ matrix. A complex number λ for which the set of linear equations

$$(A - \lambda I)\mathbf{x} = \mathbf{0} \quad (5.1)$$

has a nontrivial solution $\mathbf{x} \neq \mathbf{0}$ is called an **eigenvalue** of A ; the associated solution \mathbf{x} is called an **eigenvector** of A (corresponding to λ).

1 K 9 , D B#& 5 #&#& 9 , C @ H, #&. *
 K 6 1 1 , , 1
 1% 6 1 9 , , 6, 4
 #&>#C#&& L ,6, 6 K / , , .
 5 4 / 1% 1 E, , \$ 6 6 ,
 1 6 15 6 6 , 6 , , 6

In order to motivate the discussion that will follow, we begin with two familiar elementary examples.

In considering the rotation of a rigid body Ω , the *inertia matrix* is the 3×3 symmetric matrix

$$J = \begin{pmatrix} I & I'' & I''' \\ I & I & I'' \\ I'' & I'' & I''' \end{pmatrix}$$

whose diagonal elements are the moments of inertia about the axes,

$$I = \int (y^2 + z^2) d\Omega, \quad I'' = \int (z^2 + x^2) d\Omega, \quad I''' = \int (x^2 + y^2) d\Omega,$$

and whose off-diagonal elements are defined by the corresponding products of inertia

$$\begin{aligned} I'' &= I'' = - \int xy \, d\Omega, \\ I''' &= I''' = - \int yz \, d\Omega, \\ I' &= I' = - \int zx \, d\Omega. \end{aligned}$$

Then, the eigenvectors of the inertia matrix are the directions of the *principal axes of inertia* of the body, about which free steady rotation is possible, and the eigenvalues are the *principal moments of inertia* about these axes.

A second example, which involves matrices of any order, arises in the solution of systems of linear ordinary differential equations of the form

$$\frac{d}{dt} \mathbf{x} = \mathbf{A} \mathbf{x},$$

where \mathbf{x} is a vector of n elements, each of which is a function of the independent variable t , and \mathbf{A} is an $n \times n$ matrix whose elements are constants. If \mathbf{A} were a diagonal matrix, with diagonal elements $a_i = \dots$, $i = 1, 2, \dots, n$, the solution of this system would be straightforward, as each of the equations could be solved separately, giving

$$x_i(t) = x_i(0) \exp(-a_i t), \quad i = 1, 2, \dots, n.$$

When \mathbf{A} is not a diagonal matrix, suppose that we can find a nonsingular matrix \mathbf{M} such that

$$\mathbf{M}^{-1} \mathbf{A} \mathbf{M} = \mathbf{D},$$

where D is a diagonal matrix. Then, on letting

$$x = M^{-1} y,$$

we easily see that

$$\frac{dy}{dt} = M^{-1} A M y = D y.$$

The solution of this system of differential equations is straightforward, as we have just seen, and we then find that

$$x = (M y) = M^{-1} y(0) \exp(-\lambda t),$$

where $\lambda_j = d_{jj}$ is one of the diagonal elements of D . The numbers λ_j , $j = 1, 2, \dots, n$, are the eigenvalues of the matrix A , and the columns of M are the eigenvectors of A , so the solution of this system of differential equations requires the calculation of the eigenvalues and eigenvectors of the matrix A .

In systems of differential equations of this kind the matrix A is not necessarily symmetric. In that case, the problem is more difficult; if the eigenvalues of A are not distinct there may not exist a complete set of linearly independent eigenvectors, and then the matrix M will not exist.

In this chapter, we shall develop numerical algorithms for the solution of the algebraic eigenvalue problem (5.1), assuming throughout that A is a symmetric matrix. As has been noted above, the analogous problem for a nonsymmetric matrix is more involved, and will not be considered here.

Throughout this chapter, the set of all real-valued symmetric matrices of order n will be denoted by S_n ; thus, given a matrix $A = (a_{ij})$,

$$A \in S_n \iff A = A^T \text{ \& } a_{ij} = a_{ji}, \quad i, j = 1, 2, \dots, n.$$

We begin with a reminder of some fundamental properties.

1. If $A \in S_n$ and $B \in S_n$, then $A + B \in S_n$.
2. If $A \in S_n$ and $B \in S_n$, then $AB \in S_n$.
3. If $A \in S_n$ and $B \in S_n$, then $A - B \in S_n$.
4. If $A \in S_n$ and $B \in S_n$, then $AB - BA \in S_n$.
5. If $A \in S_n$ and $B \in S_n$, then $A^T B \in S_n$.
6. If $A \in S_n$ and $B \in S_n$, then $AB^T \in S_n$.
7. If $A \in S_n$ and $B \in S_n$, then $A^T B^T \in S_n$.
8. If $A \in S_n$ and $B \in S_n$, then $AB^T \in S_n$.
9. If $A \in S_n$ and $B \in S_n$, then $A^T B \in S_n$.
10. If $A \in S_n$ and $B \in S_n$, then $AB \in S_n$.
11. If $A \in S_n$ and $B \in S_n$, then $A^T B^T \in S_n$.
12. If $A \in S_n$ and $B \in S_n$, then $AB^T \in S_n$.
13. If $A \in S_n$ and $B \in S_n$, then $A^T B \in S_n$.
14. If $A \in S_n$ and $B \in S_n$, then $AB \in S_n$.

Theorem 5.1 *Suppose that $A \in \mathbb{R}^{n \times n}$; then, the following statements are valid.*

(i)

5.2 The characteristic polynomial

Given that A is $n \times n$ and $n \geq 4$, it is quite easy to write down the characteristic polynomial $\det(A - \lambda I)$ by expanding the determinant, and then find the roots of this polynomial of degree n in order to determine the eigenvalues of A . If $n > 4$ there is no general closed formula for the roots of a polynomial in terms of its coefficients, and therefore we have to resort to a numerical technique. A further difficulty is that the roots may be very sensitive to small changes in the coefficients of the polynomial, and we find that the effect of rounding errors in the construction of the characteristic polynomial is usually catastrophic.

Example 5.1 Consider, for example, the diagonal matrix of order 16 whose diagonal elements are $j + \frac{1}{j}$, $j = 1, 2, \dots, 16$; the eigenvalues are, of course, just the diagonal elements. Constructing the characteristic polynomial, working with 10 significant digits throughout, gives the result

$$141.3333333 \lambda^{16} + 9193.333333 \lambda^{15} + \dots$$

Using a standard numerical algorithm (such as Newton's method) for computing the roots of the polynomial and working with 10 significant digits gives the smallest root as 1.333333331, which is nearly correct to 10 significant digits. The three largest roots, however, are computed as, approximately, 15.5, 1.31 and 16.7, which are very different from their true values $14.\dot{3}$, $15.\dot{3}$, $16.\dot{3}$, respectively, even though the matrix in this example is of quite modest size, and the eigenvalues are well spaced. Thus we conclude from this example that the numerical method which constructs the characteristic polynomial and finds its roots is completely unsatisfactory for general use, except for matrices of very small size.

The fact that in general the roots of the characteristic polynomial cannot be given in closed form shows that any method must proceed by successive approximation. Although one cannot expect to produce the required eigenvalues exactly in a finite number of steps, we shall see that there exist rapidly convergent iterative methods for computing the eigenvalues and eigenvectors numerically.

5.3 Jacobi's method

This method uses a succession of orthogonal transformations to produce a sequence of matrices which approaches a diagonal matrix in the limit.

Each step in the process involves a matrix representing a plane rotation. We begin with a simple example.

Example 5.2 (The plane rotation matrix in \mathbb{R}^2) Let us suppose that $\mathbf{v} = \begin{bmatrix} x \\ y \end{bmatrix}$ and consider the matrix $\mathbf{R}(\theta) \in \mathbb{R}^{2 \times 2}$ defined by

$$\mathbf{R}(\theta) = \begin{pmatrix} \cos \theta & \sin \theta \\ -\sin \theta & \cos \theta \end{pmatrix}.$$

For a vector \mathbf{v} , $\mathbf{R}(\theta)\mathbf{v}$ is the plane rotation of \mathbf{v} around the origin by an angle θ (in the clockwise direction when $\theta > 0$ and in the anticlockwise direction when $\theta < 0$).

We note in passing that since $\cos(-\theta) = \cos \theta$, $\sin(-\theta) = -\sin \theta$ and $\cos^2 \theta + \sin^2 \theta = 1$, we have that

$$(\mathbf{R}(\theta))^{-1} = \mathbf{R}(-\theta) \quad \text{and} \quad \mathbf{R}(\theta)\mathbf{R}(-\theta) = \mathbf{I}.$$

Hence $\mathbf{R}(\theta)$ is an orthogonal matrix; i.e.,

$$\mathbf{R}(\theta)\mathbf{R}(\theta)^T = \mathbf{R}(\theta)\mathbf{R}(\theta) = \mathbf{I},$$

where \mathbf{I} is the 2×2 identity matrix.

The next definition extends the notion of plane rotation matrix to \mathbb{R}^n .

Definition 5.2 (The plane rotation matrix in \mathbb{R}^n) Suppose that $n \geq 2$, $1 \leq p < q \leq n$ and $\mathbf{v} = \begin{bmatrix} x \\ y \end{bmatrix}$. We consider the matrix $\mathbf{R}(\theta) \in \mathbb{R}^{n \times n}$ whose elements are the same as those of the identity matrix $\mathbf{I} \in \mathbb{R}^{n \times n}$, except for the four elements

$$\begin{aligned} r_{pp} &= c, & r_{qq} &= s, \\ r_{pq} &= -s, & r_{qp} &= c, \end{aligned}$$

where $c = \cos \theta$, $s = \sin \theta$.

As in Example 5.2, it is a straightforward matter to show that

$$(\mathbf{R}(\theta))^{-1} = \mathbf{R}(-\theta), \quad \mathbf{R}(\theta)\mathbf{R}(-\theta) = \mathbf{I},$$

and that, therefore,

$$\mathbf{R}(\theta)\mathbf{R}(\theta)^T = (\mathbf{R}(\theta))\mathbf{R}(\theta) = \mathbf{I}.$$

Hence $\mathbf{R}(\theta) \in \mathbb{R}^{n \times n}$ is an orthogonal matrix for any p, q such that $1 \leq p < q \leq n$, and any $\begin{bmatrix} x \\ y \end{bmatrix}$.

The basic result underlying Jacobi's method is encapsulated in the next theorem.

Theorem 5.2 Suppose that $A \in \mathbb{R}^{n \times n}$. For each pair of integers (p, q) with $1 \leq p < q \leq n$, there exists $\theta \in [0, \pi/4]$ such that the (p, q) -entry of the symmetric matrix $R = R(\theta)$ is equal to 0.

Proof For the sake of notational simplicity, we shall write R instead of $R(\theta)$ throughout the proof, and abbreviate $c = \cos \theta$ and $s = \sin \theta$.

Consider the product $A' = AR$. Evidently the only difference between A' and A is in columns p and q ; these columns of A' are linear combinations of the same two columns of A :

$$\begin{aligned} a'_i &= a_i c + a_j s \\ a'_j &= a_i s + a_j c \end{aligned}, \quad i = 1, 2, \dots, n. \tag{5.3}$$

Multiplication of A' by R on the left gives a similar result, but affects rows p and q , rather than columns p and q . Writing $B = R A'$ gives

$$\begin{aligned} b_j &= a_i c + a_s s \\ b_s &= a_i s + a_j c \end{aligned}, \quad j = 1, 2, \dots, n. \tag{5.4}$$

Combining these equations shows that $B = R A' R$, where

$$\begin{aligned} b_{ij} &= a_i c + 2a_j s c + a_i s^2, \\ b_{ji} &= a_i s + 2a_j s c + a_j c^2, \\ b_{ij} &= (a_i - a_j) s c + a_i (c^2 - s^2) = b_{ji}. \end{aligned} \tag{5.5}$$

The remaining elements of $B = R A' R$ in columns p and q are given by the expressions

$$\begin{aligned} b_{ij} &= a_i c + a_j s \\ b_{ji} &= a_i s + a_j c \end{aligned}, \quad i = 1, 2, \dots, n, \quad i \neq p, q.$$

The matrix $B = R A' R$ is evidently symmetric, so the nondiagonal elements of B in rows p and q are also given by the same expressions.

Finally, we note that all the elements of B which do not lie either in row p or q or in column p or q are the same as the corresponding elements of A , that is,

$$b_{ij} = a_{ij}, \quad \text{if } i \neq p, q \text{ and } j \neq p, q.$$

We see from (5.5) that in order to ensure that b_{pq} , the (p, q) -entry of the matrix $B = R A' R$, is equal to 0, it suffices to choose θ such that

$$\tan 2\theta = \frac{2a_{pq}}{a_{pp} - a_{qq}}; \tag{5.6}$$

thus we select

and $B = R^{-1} A R$ where R^{-1} is an orthogonal matrix, then

$$\text{Trace}(B) = \sum_{i=1}^n a_{ii} \tag{5.9}$$

The quantity

$$\sum_{i=1}^n a_{ii}^2$$

is called the **Frobenius norm** of $A \in \mathbb{R}^{n \times n}$. The Frobenius norm of $A \in \mathbb{R}^{n \times n}$ is the 2-norm of A , with A regarded as an element of a linear space of dimension n^2 over the field of real numbers; however, it is *not* a subordinate norm in the sense of Definition 2.10. In particular, the Frobenius norm on $\mathbb{R}^{n \times n}$ is not subordinate to the 2-norm on \mathbb{R}^n .

Now, one can express (5.9) equivalently by saying that the Frobenius norm of a symmetric matrix A is invariant under an orthogonal transformation: $\text{Trace}(R^{-1} A R) = \text{Trace}(A)$.

Proof of lemma The sum of squares of the elements of A is the same as the trace of A^2 , for

$$\text{Trace}(A^2) = \sum_{i=1}^n (A^2)_{ii} = \sum_{i=1}^n \sum_{j=1}^n a_{ij} a_{ji} = \sum_{i,j=1}^n a_{ij}^2 \tag{5.10}$$

since A is symmetric. Analogously, as $B = R^{-1} A R$ is symmetric, we have that

$$\text{Trace}(B^2) = \sum_{i=1}^n b_{ii}^2$$

Thus, it remains to show that $\text{Trace}(B) = \text{Trace}(A)$. Now,

$$B = (R^{-1} A R)(R^{-1} A R) = R^{-1} A^2 R, \tag{5.11}$$

since R is orthogonal. Hence B is an orthogonal transformation of A^2 which, by virtue of Theorem 5.1 (vi), means that B and A^2 have the same eigenvalues, and therefore the same trace, since the trace is the sum of the eigenvalues (see Theorem 5.1 (viii)). \square

1 9 F 9 B ? H, #. @ \$
 F 6 C * % # @ # & F 6 0 , J G 6 1 \$
 6 , 1 , 6 1 , 1 , 4 6 1 6

Now we are ready to embark on the convergence analysis of the classical Jacobi method.

Theorem 5.3 *Suppose that $A \in \mathbb{R}^{n \times n}$, $n \geq 2$. In the classical Jacobi method the off-diagonal entries in the sequence of matrices $(A^{(k)})$, generated from $A^{(0)} = A$ according to Definition 5.3, converge to 0 in the sense that*

$$\lim_{k \rightarrow \infty} \max_{i \neq j} |(A^{(k)})_{ij}| = 0. \tag{5.12}$$

Furthermore,

$$\lim_{k \rightarrow \infty} \text{Trace}(A^{(k)}) = \text{Trace}(A). \tag{5.13}$$

Proof Let a be the off-diagonal element of A with largest absolute value, and let $B = (R(\theta)) A R(\theta)$, where θ is defined by (5.7). Then, letting $c = \cos \theta$ and $s = \sin \theta$, we have that

$$\begin{pmatrix} b & b \\ b & b \end{pmatrix} = \begin{pmatrix} c & s \\ s & c \end{pmatrix} \begin{pmatrix} a & a \\ a & a \end{pmatrix} \begin{pmatrix} c & s \\ s & c \end{pmatrix},$$

and Lemma 5.1 implies that

$$b^2 + 2b^2 + b^2 = a^2 + 2a^2 + a^2.$$

Writing

$$S(A) = \sum a_{ij}^2, \quad D(A) = \sum a_{ii}^2, \quad L(A) = \sum a_{ij}^2,$$

it follows that $S(A) = D(A) + L(A)$. Now $S(B) = S(A)$ by Lemma 5.1, and so $D(B) + L(B) = D(A) + L(A)$. The diagonal entries of B are the same as those of A , except the ones in rows p and q , $1 \leq p < q \leq n$. Further, as $b = 0$, it follows that $b^2 + b^2 = a^2 + a^2 + 2a^2$. Therefore,

$$D(B) = D(A) + 2a^2.$$

Consequently,

$$L(B) = L(A) - 2a^2.$$

Now a is the largest off-diagonal element of A ; hence $L(A) \leq Na^2$ where $N = n(n-1)$ is the number of off-diagonal elements, and therefore

$$L(B) \geq (1 - 2/N)L(A). \tag{5.14}$$

On writing $A_{k+1} = A_k A_k^{-1} B$, and generating subsequent members of the sequence (A_k) in a similar manner, as indicated in the algorithm in Definition 5.3, we deduce from (5.14) that

$$|L(A_k)| \leq (1 - 2/N)^k |L(A)|, \quad k = 1, 2, 3, \dots, \quad (5.15)$$

where $N \geq 2$. Thus we conclude that $\lim_{k \rightarrow \infty} |L(A_k)| = 0$.

Now, (5.13) follows from (5.10) and (5.12) on noting that

$$\text{Trace}(A_k) = S(A_k) = S(A) = D(A_k) + L(A_k) \quad k \geq 0,$$

and passing to the limit $k \rightarrow \infty$: $\text{Trace}(A) = \lim_{k \rightarrow \infty} D(A_k)$. □

According to Theorem 5.1 (viii) the trace of A_k is the sum of the eigenvalues of A_k , and the eigenvalues of A_k are the squares of the eigenvalues of A . Thus, we have shown that the sum of the squares of the diagonal elements in the sequence of matrices (A_k) generated by the classical Jacobi method converges to the sum of the squares of the eigenvalues of A . More work is required to show that for each $i = 1, 2, \dots, n$ the sequence of diagonal elements $(a_{ii}^{(k)})$ converges to an eigenvalue of A as $k \rightarrow \infty$. We shall further discuss this question in the final paragraphs of Section 5.4. First, however, we describe another variant of Jacobi's method.

Definition 5.4 (The serial Jacobi method) *This version of Jacobi's method proceeds in a systematic order, using transformations $R_{pq}(\theta)$ to reduce to zero the elements $(1, 2), (1, 3), \dots, (1, n), (2, 3), (2, 4), \dots, (2, n), \dots, (n-1, n)$ in this order. The complete step is then repeated iteratively.*

It is not difficult to prove that this method also converges. Both these variants of the Jacobi method converge quite rapidly; the rate of convergence is in practice much faster than is suggested by (5.15), and in fact it can be shown that convergence is ultimately quadratic.

It is time for an example!

Example 5.3 *Let us consider the 5 × 5 matrix*

$$A = \begin{pmatrix} 4 & 1 & 2 & 1 & 2 \\ 1 & 3 & 0 & 3 & 4 \\ 2 & 0 & 1 & 2 & 2 \\ 1 & 3 & 2 & 4 & 1 \\ 2 & 4 & 2 & 1 & 1 \end{pmatrix}. \quad (5.16)$$

The values of $D(A^{(k)})$ and $L(A^{(k)})$ after each iteration of the serial Jacobi method, with $A^{(0)} = A$, are shown in Table 5.1. The off-diagonal elements of the third iterate, $A^{(3)}$, are zero to 10 decimal digits. The diagonal elements of $A^{(3)}$, which give the eigenvalues, are

$$8.094, 1.690, -0.671, 7.170, -3.282.$$

Note that the eigenvalues do not appear in any particular order.

Table 5.1. Convergence of the serial Jacobi iteration.

	$*$ $()_+$	$*$ $()_+$
$\&$	$\& \$$	$' \% \$) ($
$\&$	$\& \$$	$\$ ' (($
$\&$	$\&$	

This concludes the discussion about the use of Jacobi’s method for computing the eigenvalues of a symmetric matrix A . ‘Fine,’ you might say, ‘but how do we determine the *eigenvectors* of A ?’

It turns out that by collecting the information accumulated in the course of the Jacobi iteration, it is fairly easy to calculate the eigenvectors of A . We begin by noting that if M is an orthogonal matrix such that $M^{-1}AM = D$, where D is diagonal, then the diagonal elements of D are the eigenvalues of A , and the columns of M are the corresponding eigenvectors of A .

In the course of the Jacobi iteration (be it classical or serial), we have constructed the plane rotations $R_j = R(\theta_j)$, $j = 1, 2, \dots, k$. Thus, an approximation $M^{(k)}$ to the orthogonal matrix M can be obtained by considering the product of these rotation matrices: initially, we put $M^{(0)} = I$ and then we apply the column transformation $R_j = R(\theta_j)$ at each step $j = 1, 2, \dots, k$. This corresponds to multiplying $M^{(j-1)}$ on the right by $R_j = R(\theta_j)$ for $j = 1, 2, \dots, k$, and leads to the orthogonal matrix

$$M^{(k)} = R_k = R(\theta_k) \dots R_1 = R(\theta_1)$$

which represents the required approximation to the orthogonal matrix M . The columns of $M^{(k)}$ will be the desired approximate eigenvectors

of A corresponding to the approximate eigenvalues which appear along the diagonal of A .

The Jacobi method usually converges in a reasonable number of iterations, and is a satisfactory method for small or moderate-sized matrices. However, there are many problems, particularly in the area of numerical solution of partial differential equations, which give rise to very large matrices that are sparse, with most of the elements being zero. A further consideration is that in many practical situations one does not need to compute all the eigenvalues. It is much more common to require a few of the largest eigenvalues and corresponding eigenvectors, or perhaps a few of the smallest. Jacobi's method is not suitable for such problems, as it always produces all the eigenvalues, and will not preserve the sparse structure of a matrix during the course of the iteration. For example, it is easy to see that if Jacobi's method is applied to a symmetric tridiagonal matrix, then at the end of one sweep all (but two) of the elements of the matrix will in general be nonzero and, although still symmetric, the transformed matrix is no longer tridiagonal. Later on in this chapter we shall consider numerical algorithms for computing selected eigenvalues of a matrix. Thus, as an overture to what will follow, we now outline a 'rough and ready' technique for locating the eigenvalues.

5.4 The Gerschgorin theorems

Gerschgorin's Theorem provides a very simple way of determining a region that contains the eigenvalues of a matrix. It is very general, and does not assume that the matrix is symmetric; in fact we shall allow the elements of a square matrix of order n to be complex and write $A \in \mathbb{C}^{n \times n}$ to express this fact.

Definition 5.5 Suppose that $n \geq 2$ and $A \in \mathbb{C}^{n \times n}$. The **Gerschgorin discs** D_i , $i = 1, 2, \dots, n$, of the matrix A are defined as the closed circular regions

$$D_i = \{z \in \mathbb{C} : |z - a_{ii}| \leq R_i\} \quad (5.17)$$

in the complex plane, where

$$R_i = \sum_{j \neq i} |a_{ij}| \quad (5.18)$$

is the radius of D_i .

$$\begin{aligned} & \left\{ z \in \mathbb{C} : |z - a_{ii}| \leq \sum_{j \neq i} |a_{ij}| \right\} \\ & \text{where } a_{ii} \text{ is the } i\text{-th diagonal element and } \sum_{j \neq i} |a_{ij}| \text{ is the sum of the absolute values of the } i\text{-th row elements excluding } a_{ii}. \end{aligned}$$

Theorem 5.4 (Gerschgorin’s Theorem) Let $n \geq 2$ and $A \in \mathbb{C}^{n \times n}$. All eigenvalues of the matrix A lie in the region $D = \bigcup_{i=1}^n D_i$, where D_i , $i = 1, 2, \dots, n$, are the Gerschgorin discs of A defined by (5.17), (5.18).

Proof Suppose that λ and $x \neq 0$ are an eigenvalue and the corresponding eigenvector of A , so that

$$a_i x_i = \lambda x_i, \quad i = 1, 2, \dots, n. \tag{5.19}$$

Suppose that x_k , with $k = 1, 2, \dots, n$, is the component of x which has largest modulus, or one of those components if more than one have the same modulus. We note in passing that $x_k \neq 0$, given that $x \neq 0$; also,

$$x_j = \lambda^{-1} a_{jk} x_k, \quad j = 1, 2, \dots, n. \tag{5.20}$$

This means that

$$\begin{aligned} a_{jk} x_j &= \lambda^{-1} a_{jk} a_{jk} x_k \\ &= \lambda^{-1} a_{jk} a_{jk} x_k \\ &= \lambda^{-1} a_{jk} a_{jk} x_k \\ &\leq R x_k, \end{aligned} \tag{5.21}$$

which, on division by x_k , shows that λ lies in the Gerschgorin disc D_k of radius R centred at a_{kk} . Hence, $D = \bigcup_{i=1}^n D_i$. \square

Theorem 5.5 (Gerschgorin’s Second Theorem) Let $n \geq 2$. Suppose that $1 \leq p \leq n - 1$ and that the Gerschgorin discs of the matrix $A \in \mathbb{C}^{n \times n}$ can be divided into two disjoint subsets D_1 and D_2 , containing p and $q = n - p$ discs respectively. Then, the union of the discs in D_1 contains p of the eigenvalues, and the union of the discs in D_2 contains $n - p$ eigenvalues. In particular, if one disc is disjoint from all the others, it contains exactly one eigenvalue, and if all the discs are disjoint then each disc contains exactly one eigenvalue.

Proof We shall use a so-called *homotopy* (or continuation) argument.

For $0 < \alpha < 1$, we consider the matrix $B(\alpha) = (b_{ij}(\alpha))_{n \times n}$, where

$$b_{ij}(\alpha) = \begin{cases} a_{ij} & \text{if } i = j, \\ a_{ij} & \text{if } i \neq j. \end{cases} \quad (5.22)$$

Then, $B(1) = A$, and $B(0)$ is the diagonal matrix whose diagonal elements coincide with those of A . Each of the eigenvalues of $B(0)$ is therefore the centre of one of the Gerschgorin discs of A ; thus exactly p of the eigenvalues of $B(0)$ lie in the union of the discs in D_1 . Now, the eigenvalues of $B(\alpha)$ are the zeros of its characteristic polynomial, which is a polynomial whose coefficients are continuous functions of α ; hence the zeros of this polynomial are also continuous functions of α . Thus as α increases from 0 to 1 the eigenvalues of $B(\alpha)$ move along continuous paths in the complex plane, and at the same time the radii of the Gerschgorin discs increase from 0 to the radii of the Gerschgorin discs of A . Since p of the eigenvalues lie in the union of the discs in D_1 when $\alpha = 0$, and these discs are disjoint from all of the discs in D_2 , these p eigenvalues must still lie in the union of the discs in D_1 when $\alpha = 1$, and the theorem is proved.

The same proof evidently still applies when the discs can be divided into any number of disjoint subsets. \square

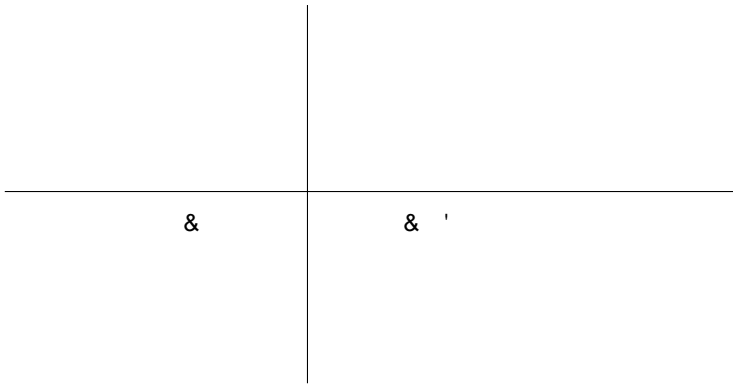
Example 5.4 Consider the matrix

$$A = \begin{pmatrix} 4.00 & 0.20 & 0.10 & 0.10 \\ 0.20 & 1.00 & 0.10 & 0.05 \\ 0.10 & 0.10 & 3.00 & 0.10 \\ 0.10 & 0.05 & 0.10 & 3.00 \end{pmatrix}. \quad (5.23)$$

Figure 5.1 shows, as solid circles, the Gerschgorin discs for this matrix; for instance, one of the discs has centre at 4.00 and radius 0.40. The discs are clearly disjoint, so that each disc contains one eigenvalue of the matrix. The significance of the dotted circles will be explained in our next example.

Example 5.5 Let us consider the matrix A defined by (5.23), and then transform it into $B = KAK^{-1}$, where $K_{n \times n}$ is the same as the identity matrix except that $k_{22} = \alpha > 0$.

This transformation has the effect of multiplying the elements in row 2 by α , and multiplying the elements in column 2 by $1/\alpha$; the diagonal element a_{22} thus remains unaltered. A small value of α then means that the second disc of B is smaller than the second disc of A , but the other



$$\begin{pmatrix} \lambda_1 & & \\ & \lambda_2 & \\ & & \dots \\ & & & \lambda_n \end{pmatrix}$$

discs grow larger. The dotted discs in Figure 5.1 are for the matrix B with $\epsilon = 1/23$. For this value the other discs are still just disjoint from the disc centred at $\lambda = 1.00$; the disc with centre at 4.00 almost touches the disc with centre at $\lambda = 1.00$. The disc with centre $\lambda = 1.00$ has radius 0.014 , and is too small to be visible in the figure. The eigenvalue in this disc is 1.009 to three decimal digits. The same procedure can be used to reduce the size of each of the discs in turn.

This idea is formalised in the next theorem.

Theorem 5.6 *Let $n \geq 2$, and suppose that in the matrix $A \in \mathbb{R}^{n \times n}$ all the off-diagonal elements are smaller in absolute value than ϵ , so that $|a_{ij}| < \epsilon$, for all $i, j = 1, 2, \dots, n$ with $i \neq j$. Suppose also that for a particular integer $r = 1, 2, \dots, n$ the diagonal element a_{rr} is distant from all the other diagonal elements, so that $a_{rr} - a_{ii} > \epsilon$, for all i such that $i \neq r$. Then, provided that*

$$\epsilon < \frac{1}{2(n-1)}, \tag{5.24}$$

there is an eigenvalue of A such that

$$|a_{rr} - \lambda| < 2(n-1) \epsilon. \quad (5.25)$$

Proof We apply the **similarity transformation**

$$A' = K^{-1} A K,$$

where K is the same as the identity matrix, except that the diagonal element in row r is chosen to be $k = \epsilon > 0$. This has the effect of multiplying the off-diagonal elements of row r by ϵ , and the element in column r of row i , where $i \neq r$, by $1/\epsilon$. The Gerschgorin disc from row r then has centre a_{rr} and radius not exceeding $(n-1)\epsilon$, and the disc corresponding to row $i \neq r$ has centre a_{ii} and radius not exceeding $(n-2)\epsilon + 1/\epsilon$.

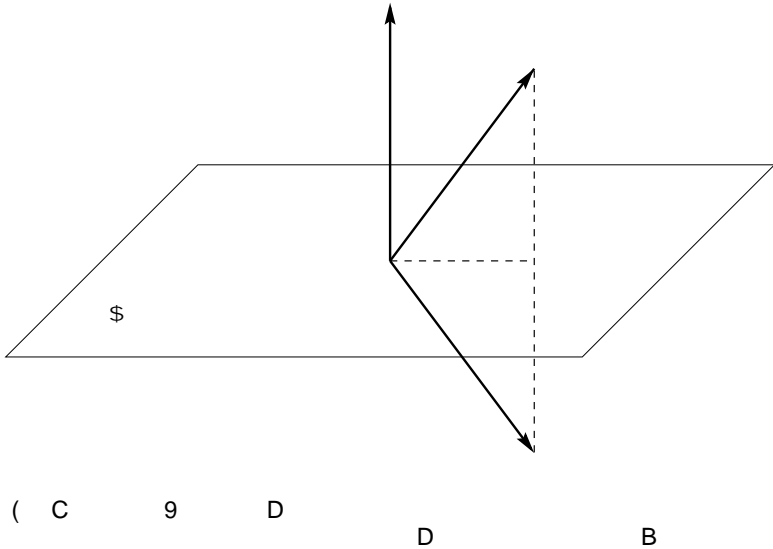
We now want to reduce the size of disc r by choosing a small value of ϵ , while keeping it disjoint from the rest. This is easily done by choosing $\epsilon = 2/(n-2)$. The radius of disc r does not exceed $2(n-1)/(n-2)$, and the radius of disc $i \neq r$ does not exceed $(n-2)/(n-2) + 1/2$. The sum of these radii therefore satisfies

$$\begin{aligned} R_r + R_i &= 2(n-1)/(n-2) + (n-2)/(n-2) + 1/2 \\ &< 2 + (n-2)/(n-2) \\ &< 3, \end{aligned} \quad (5.26)$$

where we have used the given condition (5.24) twice. As the centres a_{rr} and a_{ii} of these discs are distant more than $1/2$ from each other, (5.26) shows that the two discs are disjoint, and the required result is proved. \square

Theorem 5.6 is sufficient to show that for a matrix satisfying its hypotheses we can find a Gerschgorin disc whose radius is of order ϵ provided that ϵ is sufficiently small. It also indicates that the spacing between the diagonal elements is important.

In particular, Theorem 5.6 applies to the matrix A' which results after k iterations of the Jacobi method. If at that stage all the off-diagonal elements have magnitude less than ϵ then there is one eigenvalue in each of the intervals $[a_{ii} - (n-1)\epsilon, a_{ii} + (n-1)\epsilon]$, provided that these intervals are disjoint; this follows from Theorem 5.5. If ϵ is sufficiently small compared with the distances between the diagonal elements of A' , Theorem 5.6 may be used to give closer bounds on the eigenvalues.



hyperplane \$ consisting of all vectors that are perpendicular to in is invariant under the mapping H . Finally, for any ,

$$H = \dots$$

Hence, if the angle between and is denoted by , then the angle between and H is equal to + . We conclude from these observations that the vector H is the reflection of in the hyperplane \$. For this reason, the mapping H is frequently referred to as **Householder reflector**, corresponding to the vector (see Figure 5.2).

Lemma 5.2 *Every Householder matrix is symmetric and orthogonal.*

Proof As $I = I$, $(\dots) = (\dots) = \dots$, and is a (positive real) number, the symmetry of H follows. The orthogonality of H is a consequence of the identity

$$H H = H H = H = I - \frac{4}{(\dots)} + \frac{4}{(\dots)}(\dots)(\dots) = I,$$

since $(\dots)(\dots) = (\dots) = (\dots)$ by the associativity of matrix multiplication. □

Lemma 5.3 Let $1 \leq k < n$ and suppose that H is a $k \times k$ Householder matrix. Then, the matrix H is written in partitioned form as

$$H = \begin{pmatrix} I_{n-k} & 0 \\ 0 & H \end{pmatrix}$$

where I_{n-k} is the identity matrix of order $n - k$ and 0 is the $(n - k) \times k$ zero matrix, is also a Householder matrix.

The proof of this lemma is straightforward and is left as an exercise. (See Exercise 1.)

Lemma 5.4 Given any vector \mathbf{v} , there exists a Householder matrix H such that all elements of the vector $H\mathbf{v}$ are zero, except the first; i.e., $H\mathbf{v}$ is a nonzero multiple of \mathbf{e}_1 , the first column of the identity matrix.

In geometrical terms this result can be rephrased by saying that for any vector \mathbf{v} there exists an $(n - 1)$ -dimensional hyperplane \mathcal{H} passing through the origin in \mathbb{R}^n such that the reflection H of \mathbf{v} in \mathcal{H} is equal to a nonzero multiple of \mathbf{e}_1 . To find \mathcal{H} it suffices to identify a vector \mathbf{w} normal to \mathcal{H} . Since \mathcal{H} is unaffected by rescaling \mathbf{w} (see Definition 5.6), the length of \mathbf{w} is immaterial. As noted in the discussion following Definition 5.6, the vectors $H\mathbf{v}$, \mathbf{v} , and \mathbf{w} are coplanar. Therefore, we shall seek \mathbf{w} as a suitable linear combination of \mathbf{v} and $H\mathbf{v}$.

Proof of lemma We seek $H\mathbf{w} = \lambda \mathbf{e}_1$ with $\mathbf{w} = \mathbf{v} + cH\mathbf{v}$, where c is a nonzero real number to be determined. Hence,

$$\begin{aligned} H(\mathbf{v} + cH\mathbf{v}) &= \lambda \mathbf{e}_1 \\ H\mathbf{v} + cH^2\mathbf{v} &= \lambda \mathbf{e}_1 \end{aligned}$$

where $\lambda = \mathbf{e}_1^T \lambda \mathbf{e}_1$ is the first entry of $\lambda \mathbf{e}_1$. A simple manipulation then shows that

$$H\mathbf{v} = \frac{\lambda}{1 + 2c} (\mathbf{v} + cH\mathbf{v}) = \frac{(\mathbf{e}_1^T \lambda \mathbf{e}_1)}{1 + 2c} \mathbf{v} + \frac{2c(\mathbf{e}_1^T \lambda \mathbf{e}_1)}{1 + 2c} H\mathbf{v}.$$

Thus, $H\mathbf{v}$ will be a multiple of \mathbf{e}_1 provided that we choose c so that $1 + 2c = 0$. Also, to avoid division by 0, we need to ensure that $1 + 2c \neq 0$. To do so, note that $c \neq -1/2$; therefore

$$1 + 2c \neq 0 \implies (1 + 2c) \neq 0,$$

provided that $\alpha + c = 0$, which can be ensured by selecting the appropriate sign for c , that is, by defining

$$c = \begin{cases} (\text{sign } \alpha) \sqrt{\alpha^2 + \beta^2} & \text{when } \alpha \neq 0, \\ \beta & \text{when } \alpha = 0. \end{cases}$$

With this choice of c , we have $H = I - \frac{2}{\alpha^2 + \beta^2} \begin{pmatrix} \alpha & \beta \\ \beta & -\alpha \end{pmatrix} \begin{pmatrix} \alpha \\ \beta \end{pmatrix} \begin{pmatrix} \alpha & \beta \end{pmatrix}$, as required. □

We now show how Householder matrices can be used to reduce a given matrix to tridiagonal form.

Theorem 5.7 *Given that $A \in \mathbb{R}^{n \times n}$ and $n \geq 3$, there exists a matrix $Q \in \mathbb{R}^{n \times n}$, a product of $n - 2$ Householder matrices $H_k \in \mathbb{R}^{n \times n}$, $k = 2, \dots, n - 1$, given by*

$$Q = H_{n-2} H_{n-3} \dots H_2$$

such that $Q^T A Q = T$ is tridiagonal; the matrix Q is orthogonal.

Proof The proof of the theorem will proceed by induction. Before embarking on this, we make some preparatory observations which highlight the key ideas in the proof.

Consider the matrix $A \in \mathbb{R}^{n \times n}$, partitioned by its first row and column in the form

$$A = \begin{pmatrix} a_{11} & c^T \\ c & A_{22} \end{pmatrix},$$

where $a_{11} \in \mathbb{R}$, $c \in \mathbb{R}^{(n-1) \times 1}$ and $A_{22} \in \mathbb{R}^{(n-1) \times (n-1)}$, and define

$$\alpha = \|c\|_2, \quad \beta = \begin{cases} \alpha & \text{if } a_{11} \geq 0 \\ -\alpha & \text{if } a_{11} < 0 \end{cases} \text{ for some } \alpha \geq 0.$$

If $\alpha = 0$ happens to belong to $\mathbb{R}^{(n-1) \times 1}$, then, by Lemma 5.4, there exists an $(n-1) \times (n-1)$ Householder matrix H_{22} such that each element of $H_{22} c$, except the first, is equal to 0. If, on the other hand, $\alpha > 0$, then $H_{22} c = \alpha e_1$, trivially. Either way, $H_{22} c = \alpha e_1$.

Let us extend the Householder matrix $H_{22} \in \mathbb{R}^{(n-1) \times (n-1)}$, using Lemma 5.3 with $k = n - 1$, to a Householder matrix $H_{22} \in \mathbb{R}^{(n-1) \times (n-1)}$ by defining the $(1, 1)$ -entry of H_{22} as 1 and choosing the remaining entries in the first row and first column of H_{22} as 0. Then,

$$\begin{aligned} H_{22} A H_{22} &= \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & H_{22} \end{pmatrix} \begin{pmatrix} a_{11} & c^T \\ c & A_{22} \end{pmatrix} \begin{pmatrix} 1 & \mathbf{0} \\ \mathbf{0} & H_{22} \end{pmatrix} \\ &= \begin{pmatrix} a_{11} & \alpha e_1^T \\ \alpha e_1 & H_{22} A_{22} H_{22} \end{pmatrix} = D, \end{aligned} \tag{5.27}$$

where

=

orthogonal, $Q^T Q = I$ is itself orthogonal. Moreover, for any $f \in \mathbb{R}^n$ we have $Q f = \|f\| e_1$, since the $(1, 1)$ -entry of Q is 1 and the remaining entries in the first column of Q are 0. This concludes the inductive step, and completes the proof. \square

The recursive transformation of a symmetric matrix to tridiagonal form outlined in the proof of Theorem 5.7 is called **Householder's method**. In implementing this method in practice it is important to carry out the transformations efficiently. Counting the arithmetic operations involved is straightforward but tedious, and shows that the complete reduction requires approximately $\frac{1}{3}n^3$ multiplications, for a moderately large value of n .

Example 5.6 *In order to illustrate Householder's method, we return to the matrix A defined in (5.16). The first stage uses the Householder matrix defined by the vector*

$$v = (0.000, 4.162, 2.000, 1.000, 2.000)^T. \quad (5.28)$$

The result of the transformation is the matrix

$$\begin{pmatrix} 4.000 & 3.162 & 0.000 & 0.000 & 0.000 \\ 3.162 & 5.300 & 1.232 & 0.332 & 0.284 \\ 0.000 & 1.232 & 1.653 & 3.312 & 0.275 \\ 0.000 & 0.332 & 3.312 & 5.149 & 1.123 \\ 0.000 & 0.284 & 0.275 & 1.123 & 3.102 \end{pmatrix}.$$

The leading element of the matrix is unchanged, and the first row and column have tridiagonal structure.

The second stage uses the Householder matrix with the vector

$$v = (0.000, 0.000, 2.540, 0.332, 0.284)^T \quad (5.29)$$

and gives the new matrix

$$\begin{pmatrix} 4.000 & 3.162 & 0.000 & 0.000 & 0.000 \\ 3.162 & 5.300 & 1.308 & 0.000 & 0.000 \\ 0.000 & 1.308 & 0.057 & 2.166 & 0.792 \\ 0.000 & 0.000 & 2.166 & 6.610 & 0.420 \\ 0.000 & 0.000 & 0.792 & 0.420 & 2.967 \end{pmatrix}.$$

This time the leading 2×2 minor is unaltered, and the first two rows and columns have tridiagonal structure.

The final stage uses the Householder matrix with vector

$$v = (0.000, 0.000, 0.000, 4.471, 0.792)^T \quad (5.30)$$

The determinants of the successive principal minors of a matrix of this form can easily be calculated by recurrence. Defining $p_r(\lambda)$ to be the determinant of the leading principal minor of order r of T , we see that

$$\begin{aligned} p_1(\lambda) &= a_{11} - \lambda, \\ p_r(\lambda) &= (a_{rr} - \lambda)p_{r-1}(\lambda) - b_{r-1}a_{r-1,r}. \end{aligned}$$

Expanding $p_r(\lambda)$ in terms of the elements of the last row, and then in terms of the last column, we obtain the relation

$$p_r(\lambda) = (a_{rr} - \lambda)p_{r-1}(\lambda) - b_{r-1}p_{r-2}(\lambda), \quad r = 2, 3, \dots, n,$$

with the convention that

$$p_0(\lambda) = 1.$$

In the rest of this section we shall assume that all the off-diagonal elements b_k are nonzero. For suppose that $b_k = 0$ for some k in the set $\{2, 3, \dots, n\}$; then, the eigenvalues of the matrix T comprise the eigenvalues of the matrix consisting of the first $k-1$ rows and columns, together with the eigenvalues of the matrix consisting of the last $n-k+1$ rows and columns. These two problems become separated and can be treated independently; if several of the off-diagonal elements are zero, the matrix can be partitioned into a number of smaller matrices which can then be dealt with independently.

Theorem 5.8 (Cauchy's Interlace Theorem) *Let $n \geq 3$. The roots of p_r separate those of p_{r-1} , for $r = 1, 2, \dots, n-1$; i.e., between two consecutive roots of p_r there is exactly one root of the polynomial p_{r-1} , $r = 1, 2, \dots, n-1$.*

Proof The proof is by induction. It is trivial to show that the property holds for $r = 1$: the two roots

$$-\frac{a_{11}}{b_{12}} \text{ and } \frac{a_{11} + a_{22}}{(a_{11} - a_{22}) + 4b_{12}}$$

of p_1 are separated by a_{11} , the only root of the linear polynomial p_0 .

Suppose that the statement is true when $r = i-1$, $2 \leq i \leq n-1$, so that the roots of p_{i-1} separate those of p_{i-2} . On denoting by α and β two consecutive roots of p_{i-1} , the inductive hypothesis implies that p_{i-2} has exactly one root between α and β , which means that $p_{i-1}(\alpha)$ and

$p_{-}(\lambda)$ have opposite signs. Now,

$$p_{-}(\lambda) = (a - \lambda)p_{+}(\lambda) - b p_{-}(\lambda),$$

so that, as λ_1 and λ_2 are roots of p_{+} , it follows that $p_{-}(\lambda_1)$ and $p_{-}(\lambda_2)$ also have opposite signs. Hence p_{-} has at least one root between λ_1 and λ_2 . Choosing λ_1 and λ_2 to be each pair of consecutive roots of p_{+} in turn we have therefore located $i - 1$ roots of p_{-} .

Next choose λ_1 to be the algebraically smallest root of p_{+}

of p_{-} . Then, there is exactly one root of p between α and β ; denote this root by γ . As we saw in the proof of the previous theorem $p(\alpha)$ is positive when α is large and negative, and the sign of $p(\beta)$ is determined by the number of roots of p which are less than β . Hence if $\alpha < \beta$ both p and p_{-} have the same number of roots less than β , so that $p(\beta)$ and $p_{-}(\beta)$ have the same sign, and $s(\beta) = s_{-}(\beta) + 1$. Also if p and p_{-} have the same number of roots less than α , then p must have one more root which is greater than α ; this means that $g(\alpha) = g_{-}(\alpha) + 1$. Hence $s(\alpha) = g(\alpha)$. A similar argument shows that $s(\alpha) = g(\alpha)$ in the alternative situation where $\alpha > \beta$. It is also a simple matter to modify the argument slightly for the cases where α is less than the smallest root of p_{-} , or greater than the largest root of p_{-} , and so the inductive step is complete. \square

The theorem and proof do not allow for any of the members of the sequence being zero, in which case the sign becomes undefined. A more careful analysis is tedious but not difficult; it shows that the theorem still holds if we adopt the convention that when $p(\alpha)$ is zero it is given the same sign as $p_{-}(\alpha)$. As we have already seen, two consecutive members of the sequence cannot both be zero.

Our next example will illustrate the application of the Sturm sequence property.

Example 5.7 Determine the second largest eigenvalue of the matrix

$$A = \begin{pmatrix} 3 & 1 & 0 & 0 \\ 1 & 1 & 2 & 0 \\ 0 & 2 & 1 & 1 \\ 0 & 0 & 1 & 1 \end{pmatrix}. \tag{5.32}$$

If the eigenvalues are $\lambda_1, \lambda_2, \lambda_3, \lambda_4$, where $\lambda_1 > \lambda_2 > \lambda_3 > \lambda_4$, we wish to find λ_2 . Now, it is easy to see from Theorem 5.5 that all the eigenvalues lie in the interval $[-4, 4]$. We take the midpoint of this interval, and evaluate the Sturm sequence with $\alpha = 0$, giving

$$p(0) = 1, \quad p_{-}(0) = 3, \quad p(0) = -4, \quad p_{-}(0) = 16, \quad p(0) = 12.$$

In this sequence there are three agreements of sign:

$$(1, 3), \quad (-4, 16) \quad \text{and} \quad (16, 12).$$

Hence $s(0) = 3$, and the matrix has three eigenvalues greater than 0; this means that λ_2 must lie in the right-hand half of the interval

[4, 4], that is, in [0, 4]. We construct the Sturm sequence for $\lambda = 2$, the midpoint of the interval, giving

$$p_0(\lambda) = 1, \quad p_1(\lambda) = \lambda - 4, \quad p_2(\lambda) = 4 - \lambda, \quad p_3(\lambda) = 0, \quad p_4(\lambda) = 4.$$

Notice that here $p_2(\lambda)$ is zero, and is given the negative sign to agree with $p_1(\lambda)$. The number of agreements in sign here is two, so two of the eigenvalues are greater than 2, and $\lambda = 2$ must lie in [2, 4], the right-hand half of the interval [0, 4]. For $\lambda = 3$ we obtain the sequence

$$1, \quad +0, \quad 1, \quad 2, \quad 3,$$

with only one agreement of sign, so this time $\lambda = 3$ must lie in the left-hand half [2, 3] of the interval [2, 4], and we repeat the process, taking $\lambda = 2.5$, the midpoint of [2, 3]. This time the sequence is

$$1, \quad \frac{1}{2}, \quad \frac{11}{4}, \quad \frac{17}{8}, \quad \frac{7}{16},$$

with one agreement in sign, showing that $\lambda < 2.5$.

The process of bisection can be repeated as many times as required to locate the eigenvalue to a given accuracy. After 13 stages we find that $\lambda = 2.450$ correct to three decimal digits.

This method is very similar to the usual bisection process for finding a solution of $f(x) = 0$, beginning with an interval $[a, b]$ such that $f(a)$ and $f(b)$ have opposite signs. A great advantage of the Sturm sequence method is that it not only determines the eigenvalue, but also indicates which eigenvalue it is. If we used the Jacobi method of Section 5.3 we would have to determine *all* the eigenvalues, sort them into order, and then choose the second largest eigenvalue as λ_2 .

The Sturm sequence method will also determine how many eigenvalues of a matrix lie in a given interval (α, β) ; all that we need is to construct the Sturm sequences $(p_0(\lambda), \dots, p_n(\lambda))$ and $(p_0(\lambda), \dots, p_n(\lambda))$; then, the required number of eigenvalues is $S(\alpha) - S(\beta)$.

It is very important to calculate the sequence $p_i(\lambda)$ directly from the recurrence relation. For instance, in Example 5.7, with $\lambda = 2.445$ we obtain

$$\begin{aligned} p_0(2.445) &= 1, \\ p_1(2.445) &= 3 - 2.445 = 0.555, \\ p_2(2.445) &= (1 - 2.445) \cdot 0.555 - 1 = -2.9120, \\ p_3(2.445) &= (1 - 2.445) \cdot 2.9120 - 4 \cdot 0.555 = 1.9878, \\ p_4(2.445) &= (1 - 2.445) \cdot 1.9878 - 1 \cdot 2.9120 = 0.0396. \end{aligned}$$

The alternative, to construct explicit forms for the polynomials $p_j(\lambda)$, $j = 0, 1, \dots, n$, and then evaluate $p_j(\lambda)$ by inserting the value of $\lambda = \lambda_k$ into each of the polynomials $p_j(\lambda)$, will lead to the construction of the explicit form of the characteristic polynomial of the matrix, which is $p_n(\lambda)$, and we have already seen that this is affected disastrously by rounding errors. The calculation by direct use of the recurrence relation is perfectly satisfactory.

Example 5.8 As a second example, we return to the matrix A in (5.16), which has been transformed to the tridiagonal form (5.31), to determine the largest eigenvalue.

Table 5.2. Bisection process for the largest eigenvalue. In the table k denotes the iteration number, λ_k the k th iterate approximating the unknown eigenvalue λ , and $S(\lambda_k)$ signifies the number of sign agreements in the Sturm sequence $p_0(\lambda_k), \dots, p_n(\lambda_k)$.

k	λ_k	$S(\lambda_k)$
1	8.000000	1
2	5.000000	5
3	6.000000	8
4	6.500000	6
5	6.750000	5
6	6.875000	4
7	6.937500	3
8	6.968750	2
9	6.984375	1
10	6.992187	0

Table 5.2 shows the result of the bisection process, using the Sturm sequence. The ∞ -norm of the tridiagonal matrix is 10.926, so the process begins with the interval $[0, 10.926]$. The largest eigenvalue

is 8.094, to three decimal digits, agreeing with the result of Jacobi's method, in Section 5.3. This table also shows how some savings are possible when all the eigenvalues are required. We see from the table that use of $\epsilon = 7.511$ gives 1 agreement in sign, while $\epsilon = 6.829$ gives 2 agreements in sign. The bisection process for the second largest eigenvalue can therefore begin with the interval $[6.829, 7.511]$.

The method of bisection may appear rather crude, but it has the great advantage of guaranteed success, and is very little affected by rounding errors. Moreover, the amount of work involved is not large. If we have calculated the squares of the off-diagonal entries, b_{ij} , of the matrix T in advance, each computation of all members of the sequence requires about $2n$ multiplications. If the bisection process is continued for 40 stages, the eigenvalue will be determined to about nine significant digits, and if we require to calculate m of the eigenvalues to this accuracy, we shall need about $80mn$ multiplications. If m is a good deal smaller than n , the order of the matrix, this is likely to be a great deal smaller than the work involved in the process of reduction to tridiagonal form, which, as we have seen, is about $\frac{1}{2}n^2$ multiplications. In most practical problems it is the initial Householder reduction to tridiagonal form which accounts for most of the computational work.

5.7 The QR algorithm

In this section we discuss briefly the QR algorithm, an alternative method for determining the eigenvalues of a tridiagonal matrix. In principle it could be applied to a full matrix, but it is more efficient to use the Householder method to reduce the matrix to tridiagonal form first. The basis of the method is the QR factorisation of the matrix which we have already encountered in Chapter 2, in the solution of least squares problems. In contrast with Section 2.9, however, where we were concerned with the solution of least squares problems for rectangular matrices $A \in \mathbb{R}^{m \times n}$, here the focus is on eigenvalue problems for symmetric tridiagonal matrices $A \in \mathbb{R}^{n \times n}$; we shall therefore revisit the derivation of the QR factorisation by adopting a slightly different approach from the one proposed in Section 2.9.

5.7.1 The QR factorisation revisited

Suppose that $n \geq 3$ and $A \in \mathbb{R}^{n \times n}$ is a symmetric tridiagonal matrix. We first show how to construct an orthogonal matrix $Q \in \mathbb{R}^{n \times n}$ and

an upper triangular matrix \mathbf{R} \times such that $\mathbf{A} =$

which are already equal to zero, remain zero upon multiplication by the next rotation matrix Q in the sequence, we deduce that, after successive multiplications of A on the left by Q_1, Q_2, \dots, Q_{n-1} , the matrix

$$Q_{n-1} Q_{n-2} \dots Q_1 A = R, \tag{5.35}$$

is upper triangular. In fact, since A is tridiagonal, R is tridiagonal and upper triangular; consequently, R is **bidiagonal** in the sense that $R_{ij} = 0$ if $i = j + 1$.

As the matrices $Q_p = R_{p+1, p+1}^{-1} \begin{pmatrix} \cos \theta_p & \sin \theta_p \\ -\sin \theta_p & \cos \theta_p \end{pmatrix}$, $p = 1, 2, \dots, n - 1$, are orthogonal, and therefore $Q_p Q_p^T = I$, on multiplying (5.35) on the left by $Q_1^T Q_2^T \dots Q_{n-1}^T$, we find that

$$A = QR,$$

where

$$Q = Q_1^T Q_2^T \dots Q_{n-1}^T$$

is an orthogonal matrix (as it is a product of orthogonal matrices). The next subsection describes the QR algorithm, based on the QR factorisation, for the numerical solution of the eigenvalue problem (5.1) where the matrix $A \in \mathbb{R}^{n \times n}$ is symmetric and tridiagonal.

5.7.2 The definition of the QR algorithm

Suppose that $A \in \mathbb{R}^{n \times n}$ is symmetric and tridiagonal. The QR algorithm defines a sequence of symmetric tridiagonal matrices $A_k \in \mathbb{R}^{n \times n}$, $k = 0, 1, 2, \dots$, starting with $A_0 = A$, as follows.

Suppose that $k \geq 0$. The k th step of the QR algorithm takes the symmetric tridiagonal matrix A_k and chooses a **shift** μ_k (the choice of μ_k will be discussed below), then forming the QR factorisation

$$A_k - \mu_k I = Q_k R_k.$$

We then multiply Q_k and R_k in the reverse order, and construct the new matrix A_{k+1} defined by

$$A_{k+1} = R_k Q_k + \mu_k I.$$

Recalling that the matrix Q_k is orthogonal, it is a simple matter to see that $A_{k+1} = Q_k A_k Q_k^T$, so that A_{k+1} and A_k have the same eigenvalues. As $A_0 = A$, all matrices in the sequence (A_k) have the same eigenvalues as A itself. It is also easy to show that each of the matrices A_k is symmetric and tridiagonal. (See Exercise 7.)

The choice of the shift parameter μ is very important; if correctly chosen the sequence of matrices A_k converges very rapidly to a matrix in which one of the off-diagonal elements is zero. If this element is in the first or last row, we have thereby identified one of the eigenvalues; if it is one of the intermediate elements, we can split the matrix into two separate matrices of lower order. In either case we can repeat the iterative process with smaller matrices, until all the eigenvalues are found.

The usual simple choice of the shift parameter in the k th step is

$$\mu = a_{nn}$$

the last diagonal element of the matrix A_k . In general, after a few steps of the iteration the element at position $(n, n-1)$ will become negligibly small. One of the eigenvalues of the resulting matrix is then the last diagonal element, and we continue the process with the matrix of order $n-1$ obtained by removing the last row and column. There are special circumstances where this choice of shift is unsatisfactory, and other situations where another choice is more efficient, but we shall not discuss the details any further. The proof of the convergence of this method is long and technical; details will be found in the books cited in the Notes at the end of the chapter.

The method does not determine the eigenvalues in any particular order, so if we require only a small number of the largest eigenvalues, for example, the Sturm sequence method is preferable. The usual recommendation is that the QR algorithm should be used on a matrix of order n if more than about $\frac{1}{2}n$ of the eigenvalues are required.

Example 5.9 We apply the QR algorithm to the tridiagonal matrix (5.31).

After one step of the iteration the matrix $A_1 = R_1 Q_1 + \mu I$, with $\mu = a_{nn} = a_{55}$, is

$$A_1 = \begin{pmatrix} 7.034 & 2.271 & 0 & 0 & 0 \\ 2.271 & 2.707 & 0.744 & 0 & 0 \\ 0 & 0.744 & 5.804 & 3.202 & 0 \\ 0 & 0 & 3.202 & 0.464 & 1.419 \\ 0 & 0 & 0 & 1.419 & 2.082 \end{pmatrix}.$$

In successive iterations $k = 1, 2, 3, 4, 5$, the element a_{54} has the values 1.419, 1.262, 0.965, 0.223, 0.002, and after the next iteration a

vanishes to 10 decimal digits. The element a_{11} is 3.282, which is therefore an eigenvalue.

We then remove the last row and column, and continue the process on the resulting 4×4 matrix. After just one iteration the element at position $(4, 3)$ vanishes to 7 decimal digits, giving the eigenvalue 0.671. We remove the last row and column and continue with the resulting 3×3 matrix. After one iteration of the resulting 3×3 matrix the element at position $(3, 2)$ is 0.0005, and another iteration gives the accurate eigenvalue 1.690. We are now left with a 2×2 matrix, and the calculation of the last two eigenvalues is trivial. The number of iterations required to isolate each eigenvalue reduces as the algorithm reduces the size of the matrix; this sort of behaviour is typical.

The numerical values agree with those obtained by Jacobi's method, and the bisection method.

5.8 Inverse iteration for the eigenvectors

We saw in Section 5.3 that Jacobi's method can also, if required, produce the eigenvectors of the matrix, but the use of Householder's algorithm, in conjunction with the Sturm sequence method or the QR algorithm, only gives the eigenvalues. Suppose that $A \in \mathbb{R}^{n \times n}$ is a symmetric matrix, and assume that we have a good approximation λ_k to the required eigenvalue λ_k of A , and some approximation \mathbf{v}_k , $\|\mathbf{v}_k\| = 1$, to the associated eigenvector \mathbf{v}_k , $\|\mathbf{v}_k\| = 1$. It is implicitly assumed that $\lambda_k \neq \lambda_j$ and that λ_k is not an eigenvalue of A , so that the matrix $A - \lambda_k I$ is nonsingular. The method of **inverse iteration** defines the sequence of vectors \mathbf{w}_k , $k = 0, 1, \dots$, as follows: given \mathbf{w}_0 , find \mathbf{w}_1 and then \mathbf{w}_2 from

$$\begin{aligned} (A - \lambda_k I) \mathbf{w}_k &= \mathbf{w}_{k-1}, \\ \mathbf{w}_k &= \mathbf{c}_k \mathbf{w}_{k-1}, \end{aligned} \tag{5.36}$$

where $\mathbf{c}_k = 1 / \|\mathbf{w}_{k-1}\| = 1 / \|\mathbf{w}_{k-1}\|$. Hence, we conclude that $\mathbf{w}_k = \mathbf{c}_k \mathbf{w}_{k-1}$, $k = 0, 1, 2, \dots$

Theorem 5.10 *Suppose that $A \in \mathbb{R}^{n \times n}$. The sequence of vectors (\mathbf{w}_k) in (5.36) defined in the process of inverse iteration (5.36) converges to the normalised eigenvector \mathbf{v}_k corresponding to the eigenvalue λ_k which is closest to λ_k , provided that λ_k is a simple eigenvalue and the initial vector \mathbf{w}_0 is not orthogonal to the vector \mathbf{v}_k .*

Proof According to Theorem 5.1 (vii), the vector \mathbf{v}_k can be expressed

If the estimate λ_j is within rounding error of λ and the eigenvalues are well spaced, the convergence of the sequence (5.36) will be extremely rapid: usually a couple of iterations will be sufficient.

The proof of Theorem 5.10 breaks down if $\langle x, v_j \rangle = 0$, *i.e.*, when the initial vector x is exactly orthogonal to the required eigenvector. However, this does not mean that the iteration (5.36) will also break down; for the effect of rounding error will almost always introduce a small multiple of the vector v_j into the expansion of x in terms of the v_j with $j = 1, 2, \dots, n$, and the required eigenvector will then be obtained in a small number of iterations. This is a useful property of the method, since in practice it is not possible to check whether or not x is orthogonal to v_j , given that the eigenvector v_j is unknown.

There will also be a problem if there is a multiple eigenvalue, or two eigenvalues are very close together: in the first case $\lambda_j / \lambda_{j+1} = 1$ for some $j = s$, and the proof of Theorem 5.10 breaks down; in the second case $|\lambda_j - \lambda_{j+1}| / \lambda_j \approx 1$ for some $j = s$, leading to very slow convergence.

The computation of x from (5.36) requires the solution of a system of linear equations whose matrix is $A - \lambda_j I$. This matrix will usually be nearly singular – in fact, our objective in choosing λ_j was to make $A - \lambda_j I$ exactly singular. In general the solution of such a system is extremely dangerous, because of the effect of rounding errors; in this case, however, the effect of rounding error will be to introduce a multiple of the dominant eigenvector, and this is exactly what is required. An analysis of the effect of rounding errors will confirm this fact, but would take too long here.

There are two ways in which we can implement the inverse iteration process. One obvious possibility would be to use the original matrix $A - \lambda_j I$, as implied in (5.36). An alternative is to replace A in this equation by the tridiagonal matrix $T - \lambda_j I$ supplied by Householder’s method. The calculation is then very much quicker, but produces the eigenvector of T ; to obtain the corresponding eigenvector of A we must then apply to this vector the sequence of Householder transformations which were used in the original reduction to tridiagonal form. It is easy to show that this is the most efficient method.

1 9 1 1 1 4 4* 4
 , CO J K #.
 F 494 (: % #
 5 6 #??4

Inverse iteration with the original matrix A requires the LU decomposition of A , followed by one or more forward and backsubstitution operations. As we saw in Section 2.6, the LU decomposition requires approximately n^2 multiplications. The same process with the tridiagonal matrix T , using the Thomas algorithm, involves only a small multiple of n multiplications.

Having found an eigenvector of the tridiagonal matrix T , so that

$$T \mathbf{q} = \lambda \mathbf{q},$$

we use the fact that $Q^{-1}AQ = T$ to write

$$AQ = Q \mathbf{q},$$

so that the vector \mathbf{q} is an eigenvector of A . Using Theorem 5.7, this means that the required eigenvector of A is

$$\mathbf{H}_1 \mathbf{H}_2 \dots \mathbf{H}_{n-1} \mathbf{q},$$

where the matrices \mathbf{H}_j , $j = 2, \dots, n-1$, are Householder matrices. To multiply a vector \mathbf{v} by a Householder matrix $\mathbf{H} = \mathbf{H}(\mathbf{u})$ we write

$$\mathbf{H} \mathbf{v} = (\mathbf{I} - 2\mathbf{u}\mathbf{u}^T) \mathbf{v} = \mathbf{v} - 2(\mathbf{u}^T \mathbf{v}) \mathbf{u}.$$

Assuming that $\mathbf{u}^T \mathbf{v} = 2/(\mathbf{u}^T \mathbf{u})$ is known, this requires the calculation of the scalar product $\mathbf{u}^T \mathbf{v}$, and then subtracting a multiple of the vector \mathbf{u} from the vector \mathbf{v} . This evidently involves $2n$ multiplications. Hence the calculation of \mathbf{q} requires only $2n(n-2)$ multiplications, and the work involved in the whole process is proportional to n^2 , instead of n^3 . In fact the total is less than $2n(n-2)$, since a more careful count can use the fact that many of the elements in the vector \mathbf{q} are known to be zero.

Example 5.10 Returning to the tridiagonal matrix (5.31), the QR algorithm has given an accurate eigenvalue which is 8.094 to three decimal digits. Beginning the inverse iteration (5.36) with a randomly chosen vector $\mathbf{q} = (1, 0, 0, 0, 0)^T$, we find that

$$\mathbf{q} = (0.0249, 0.0574, 0.3164, 0.4256, 0.8455)^T.$$

Successive iterations make no change in this vector, as might be expected, since the eigenvalue used was accurate to within rounding error.

This is therefore the eigenvector of the tridiagonal matrix (5.31), to

eigenvectors \mathbf{v}_j , $j = 1, 2, \dots, n$, as

$$\mathbf{A} \mathbf{v}_j = \lambda_j \mathbf{v}_j, \tag{5.41}$$

then

$$R(\mathbf{v}_j) = \lambda_j. \tag{5.42}$$

On noting that δ_{ij} is equal to 1 when $i = j$ and to 0 otherwise, (5.42) follows trivially by inserting (5.41) into (5.40).

Theorem 5.12 Let $\mathbf{A} \in \mathbb{R}^{n \times n}$. For any vector \mathbf{v} ,

$$R(\mathbf{v}) \in [\lambda_{\min}, \lambda_{\max}], \tag{5.43}$$

where λ_{\min} and λ_{\max} are respectively the least and greatest of the eigenvalues of \mathbf{A} . These bounds are attained when \mathbf{v} is the corresponding eigenvector.

Proof The inequalities follow immediately from (5.42) by noting that $\delta_{ij} \geq 0$, $j = 1, 2, \dots, n$. □

Theorem 5.13 Suppose that \mathbf{v}_k is a normalised vector, that is, $\|\mathbf{v}_k\| = 1$. Assume, further, that \mathbf{v}_k is the k th normalised eigenvector of $\mathbf{A} \in \mathbb{R}^{n \times n}$, and that

$$\mathbf{A} \mathbf{v}_k = \lambda_k \mathbf{v}_k$$

for a small ϵ . Then,

$$R(\mathbf{v}_k) = \lambda_k + \mathcal{O}(\epsilon^2).$$

Proof It follows from (5.41) that $\mathbf{A} \mathbf{v}_k = \lambda_k \mathbf{v}_k$, and therefore,

$$\begin{aligned} R(\mathbf{v}_k) &= \frac{\mathbf{v}_k^T \mathbf{A} \mathbf{v}_k}{\mathbf{v}_k^T \mathbf{v}_k} \\ &= \frac{\mathbf{v}_k^T (\lambda_k \mathbf{v}_k)}{\mathbf{v}_k^T \mathbf{v}_k} \\ &= \lambda_k \frac{\mathbf{v}_k^T \mathbf{v}_k}{\mathbf{v}_k^T \mathbf{v}_k} = \lambda_k. \end{aligned}$$

Hence, $R(\mathbf{v}_k) = \lambda_k$. Further,

$$1 = \mathbf{v}_k^T \mathbf{v}_k = \mathbf{v}_k^T \mathbf{A} \mathbf{v}_k / \lambda_k$$

$$\begin{aligned}
 &= \quad + \\
 &= 1 + (\quad) + \quad .
 \end{aligned}$$

Consequently, $\quad = (\quad)$ for all $j = k$. The result then follows from (5.42) which (with $\quad = \quad = 1$) yields that

$$\begin{aligned}
 R(\quad) &= \quad + \\
 &= \quad + (\quad) .
 \end{aligned}$$

□

This important result means that if we have a fairly close approximation \quad to an eigenvector of \mathbf{A} , then the Rayleigh quotient $R(\quad)$ gives very easily a much more accurate approximation to the corresponding eigenvalue.

5.10 Perturbation analysis

It is often necessary to have an estimate of how much the eigenvalues and eigenvectors of a matrix are affected by changes in the elements. Such perturbations may arise, for example, when the matrix elements are obtained by physical measurements which are inexact, or they might result from finite difference approximations of a differential equation, as will be seen in Chapter 13. The last two theorems in this chapter address some of these questions. We begin with the following preliminary result.

Theorem 5.14 *Let $\mathbf{M} \in \mathbb{R}^{n \times n}$, with eigenvalues $\lambda_1, \dots, \lambda_n$ and corresponding orthonormal eigenvectors $\mathbf{v}_1, \dots, \mathbf{v}_n$, $i = 1, 2, \dots, n$, and suppose that $\mathbf{p} = \mathbf{0}$ and \mathbf{q} are vectors in \mathbb{R}^n and μ is a real number such that*

$$(\mathbf{M} - \mu \mathbf{I}) \mathbf{p} = \mathbf{q} . \tag{5.44}$$

Then, at least one eigenvalue λ_k of \mathbf{M} satisfies

$$\mu - \lambda_k \leq \frac{\|\mathbf{q}\|}{\|\mathbf{p}\|} .$$

Proof If μ is equal to one of the eigenvalues the proof is trivial, so we shall assume that $\mu \neq \lambda_k$, $k = 1, 2, \dots, n$. We write the vectors \mathbf{p} and

as linear combinations of the eigenvectors of M , so that

$$= \sum_{k=1}^n \alpha_k v_k, \quad = \sum_{k=1}^n \beta_k v_k.$$

Substituting in (5.44), we may equate coefficients of the linearly independent vectors $v_k, k = 1, 2, \dots, n$, to deduce that

$$(\mu - \lambda_k) \alpha_k = \beta_k, \quad k = 1, 2, \dots, n.$$

Now suppose that μ is the eigenvalue which is closest to μ ; this means that

$$|\mu - \lambda_k| > |\mu - \lambda_j|, \quad k = 1, 2, \dots, n.$$

Since the eigenvectors $v_i, i = 1, 2, \dots, n$, are orthonormal in \mathbb{R}^n , we have

$$\sum_{k=1}^n \alpha_k^2 = 1, \quad \sum_{k=1}^n \beta_k^2 = |\mu - \lambda_j|^2.$$

Hence

$$\frac{|\mu - \lambda_j|^2}{(\mu - \lambda_j)^2} = \sum_{k=1}^n \left(\frac{\beta_k}{\mu - \lambda_k} \right)^2,$$

which gives

$$= \sum_{k=1}^n \left(\frac{\beta_k}{\mu - \lambda_k} \right)^2 \frac{(\mu - \lambda_j)^2}{(\mu - \lambda_k)^2} = \left(\frac{\beta_j}{\mu - \lambda_j} \right)^2,$$

as required. □

We shall now use this result to show that in the case of a symmetric matrix A , small symmetric perturbations of A lead to small changes in the eigenvalues of A .

Theorem 5.15 (Bauer–Fike Theorem (symmetric case)) *Suppose that $A, E \in \mathbb{R}^{n \times n}$ and $B = A + E$. Assume, further, that the eigenvalues of A are denoted by $\lambda_j, j = 1, 2, \dots, n$, and μ is an eigenvalue of B . Then, at least one eigenvalue λ_j of A satisfies*

$$|\mu - \lambda_j| \leq \|E\|.$$

Proof This is a straightforward consequence of the previous theorem. Suppose that v is the normalised eigenvector of B corresponding to the eigenvalue μ , so that $Bv = \mu v$. Then,

$$(A + E)v = \mu v \implies (A - \mu I)v = -Ev.$$

It then follows from Theorem 5.14 that there is an eigenvalue μ of A such that

$$\mu \mathbf{E} - \mathbf{A} = \mathbf{E} \mathbf{E}^T,$$

as required. □

Example 5.11 Consider the 3 × 3 Hilbert matrix

$$\mathbf{A} = \begin{pmatrix} 1 & 1/2 & 1/3 \\ 1/2 & 1/3 & 1/4 \\ 1/3 & 1/4 & 1/5 \end{pmatrix}$$

and its perturbation

$$\mathbf{B} = \begin{pmatrix} 1.0000 & 0.5000 & 0.3333 \\ 0.5000 & 0.3333 & 0.2500 \\ 0.3333 & 0.2500 & 0.2000 \end{pmatrix}$$

which results by rounding each entry of A to four decimal digits.

In this case, $\mathbf{E} = \mathbf{A} - \mathbf{B}$ and $\|\mathbf{E}\| = 3.3 \times 10^{-4}$. Let μ be an eigenvalue of \mathbf{B} ; then, according to Theorem 5.15, at least one of the eigenvalues $\lambda_1, \lambda_2, \lambda_3$ of the matrix \mathbf{A} satisfies the inequality

$$|\mu - \lambda_i| \leq 3.3 \times 10^{-4}. \tag{5.45}$$

Indeed, the true eigenvalues of \mathbf{A} and \mathbf{B} are, respectively,

$$\lambda_1 = 0.002687338072, \quad \lambda_2 = 0.1223270673, \quad \lambda_3 = 1.408318925,$$

and

$$\mu_1 = 0.002664493933, \quad \mu_2 = 0.1223414532, \quad \mu_3 = 1.408294053.$$

Therefore,

$$|\mu_1 - \lambda_1| = 2.29 \times 10^{-5}, \quad |\mu_2 - \lambda_2| = 1.44 \times 10^{-5}, \quad |\mu_3 - \lambda_3| = 2.49 \times 10^{-5},$$

which is in agreement with (5.45).

5.11 Notes

Theorem 5.15 is a special case of the following general result, known as the Bauer–Fike Theorem.

¹ 94: 4 4 4 9 " 8, (' #*&C # # #@? 4

Theorem 5.16 Assume that $A \in \mathbb{C}^{n \times n}$ is diagonalisable; i.e., there exists a nonsingular matrix $X \in \mathbb{C}^{n \times n}$ such that $X^{-1}AX = \Lambda$, where Λ is a diagonal matrix whose diagonal entries $\lambda_j, j = 1, \dots, n$, are the eigenvalues of A . Suppose further that $E \in \mathbb{C}^{n \times n}, B = A + E$, and μ is an eigenvalue of B . Then, at least one eigenvalue λ_j of A satisfies

$$|\mu - \lambda_j| \leq \kappa(X) \|E\|_2,$$

where $\kappa(X) = \|X\|_2 \|X^{-1}\|_2$ is the condition number of the matrix X in the matrix 2-norm on \mathbb{C}^n .

In the special case when $A, E \in \mathbb{R}^{n \times n}$, the matrix X can be chosen to be orthogonal; i.e., $X^{-1} = X^T$. Therefore, $\|X\|_2 = \|X^{-1}\|_2 = 1$, and hence $\kappa(X) = 1$, in accordance with the inequality stated in Theorem 5.15. Theorems 5.15 and 5.16 estimate how far the eigenvalues of A are perturbed by changes in the elements of A . The question as to how large the changes in the eigenvectors may be is more difficult; it is discussed in detail in

W. T. S. Arnold, *The Algebraic Eigenvalue Problem*, Clarendon Press, Oxford University Press, New York, 1988.

Chapter 8 of Wilkinson's book outlines the convergence proof of the QR iteration, while the convergence of Jacobi's method is covered in Chapter 5 of that book. For further details, see also Chapter 9 of

W. T. S. Arnold, *The Symmetric Eigenvalue Problem*, Prentice-Hall, Englewood Cliffs, NJ, 1980.

Exercises

- 5.1 Give a proof of Lemma 5.3.
- 5.2 Use Householder matrices to transform the matrix

$$A = \begin{pmatrix} 2 & 1 & 2 & 2 \\ 1 & 7 & 6 & 5 \\ 2 & 6 & 2 & 5 \\ 2 & 5 & 5 & 1 \end{pmatrix}$$

to tridiagonal form.

5.3 Use Sturm sequences to show that no eigenvalue of the matrix

$$A = \begin{pmatrix} 3 & 1 & 0 & 0 \\ 1 & 2 & 2 & 0 \\ 0 & 2 & 4 & \\ 0 & 0 & & 1 \end{pmatrix}$$

lies in the interval $(0, 1)$ if $5 > 8$, and that exactly one eigenvalue of A lies in this interval if $5 < 8$.

5.4 Given any two nonzero vectors \mathbf{u} and \mathbf{v} in \mathbb{R}^n , construct a Householder matrix H such that $H\mathbf{u}$ is a scalar multiple of \mathbf{v} ; note that if $H\mathbf{u} = c\mathbf{v}$, then $c = \|\mathbf{u}\| / \|\mathbf{v}\|$. Is the matrix unique?

5.5 Suppose that the matrix $D \in \mathbb{R}^{n \times n}$ is diagonal with distinct diagonal elements d_1, \dots, d_n . Let $A \in \mathbb{R}^{n \times n}$, with $a_{ij} = 1$ for all $i, j = 1, 2, \dots, n$, and assume that ϵ is so small that ϵ can be neglected, and that the matrix $D + \epsilon A$ has eigenvalue $d_j + \mu$ and corresponding eigenvector $\mathbf{e} + \epsilon \mathbf{u}$. Show that $\mu = d_j$ for some $j = 1, 2, \dots, n$ and that $\mu = a_{jj}$. Write down the elements of \mathbf{e} , and show that

$$\mathbf{u} = \frac{a_{ij}}{d_i - d_j}, \quad i \neq j.$$

Explain why the requirement that eigenvectors should be normalised implies that $\mathbf{u} = \mathbf{0}$.

5.6 With the same notation as in Exercise 5, suppose now that $d_1 = d_2 = \dots = d_k$, that d_1, d_{k+1}, \dots, d_n are distinct, and that ϵ can be neglected. Writing the matrices and the eigenvector in partitioned form, so that

$$\begin{pmatrix} d_1 I_k + A_{11} & A_{12} \\ A_{21} & D_{22} + A_{22} \end{pmatrix} \begin{pmatrix} \mathbf{e}_1 + \epsilon \mathbf{u}_1 \\ \mathbf{f}_1 + \epsilon \mathbf{v}_1 \end{pmatrix} = (d_1 + \mu + \epsilon \mu) \begin{pmatrix} \mathbf{e}_1 + \epsilon \mathbf{u}_1 \\ \mathbf{f}_1 + \epsilon \mathbf{v}_1 \end{pmatrix},$$

show that $\mu = d_1$, $\mathbf{f} = \mathbf{0}$, and that μ is an eigenvalue of A_{11} with corresponding eigenvector \mathbf{e}_1 . Show how \mathbf{u}_1 is obtained from the solution of $(D_{11} - d_1 I_k) \mathbf{u}_1 = A_{12} \mathbf{v}_1$, and that

$$(A_{11} - \mu I_k) \mathbf{u}_1 = \mathbf{e}_1 - A_{12} \mathbf{v}_1.$$

Explain how the vector \mathbf{v} can be obtained in terms of the eigenvectors and eigenvalues of the matrix \mathbf{A} , assuming that these eigenvalues are distinct.

- 5.7 Suppose that $\mathbf{A} \in \mathbb{R}^{n \times n}$ is tridiagonal, that $\mathbf{A} - \mu \mathbf{I} = \mathbf{QR}$ and $\mathbf{B} = \mathbf{RQ} + \mu \mathbf{I}$, where $\mu \in \mathbb{R}$, $\mathbf{Q} \in \mathbb{R}^{n \times n}$ is a product of plane rotations and $\mathbf{R} \in \mathbb{R}^{n \times n}$ is upper triangular and tridiagonal. Show that \mathbf{B} can be written as an orthogonal transformation of \mathbf{A} , and that \mathbf{B} is symmetric. Show also that the only nonzero elements in the matrix \mathbf{B} which are below the diagonal lie immediately below the diagonal; deduce that \mathbf{B} is tridiagonal.
- 5.8 Perform one step of the QR algorithm, using the shift $\mu = a_{11}$, for the matrix

$$\mathbf{A} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}.$$

Show that the QR algorithm does not converge for this matrix. (This is a special case in which a different shift must be used.)

- 5.9 Perform one step of the QR algorithm, using the shift $\mu = a_{11}$, for the matrix

$$\mathbf{A} = \begin{pmatrix} 13 & 4 \\ 4 & 10 \end{pmatrix}.$$

- 5.10 Carry out two steps of inverse iteration for the matrix

$$\mathbf{A} = \begin{pmatrix} 2 & 2 \\ 2 & 5 \end{pmatrix},$$

using the eigenvalue estimate $\lambda = 5$ and the initial vector

$$\mathbf{v} = \begin{pmatrix} 1 \\ 1 \end{pmatrix}.$$

Verify that the elements of the vector \mathbf{v} agree with those of the true eigenvector with an accuracy of about 5%. Evaluate the Rayleigh quotient using the vector \mathbf{v} , and verify that the result agrees with the true eigenvalue to about 1 in 3000.

- 5.11 An eigenvalue and eigenvector of the matrix \mathbf{A} may be evaluated by solving the system of nonlinear equations

$$\begin{aligned} (\mathbf{A} - \lambda \mathbf{I}) \mathbf{v} &= \mathbf{0}, \\ \mathbf{v}^T \mathbf{v} &= 1 \end{aligned}$$

for the unknowns λ and \mathbf{v} . Using Newton's method, starting

from estimates \bar{a} and \bar{b} , show that the next iteration is determined by

$$\begin{aligned} A_{i+1} &= (A_{ii} - \bar{a}) - \bar{b}^2 / (A_{ii} - \bar{a}), \\ &= -(\bar{b}^2 / (A_{ii} - \bar{a}) + \bar{a}) \end{aligned}$$

and $\bar{a}_{i+1} = \bar{a} + \bar{b}^2 / (A_{ii} - \bar{a})$, $\bar{b}_{i+1} = \bar{b}$. Comment on the difference between this method and the method of inverse iteration in Section 5.8.

- 5.12 Suppose that $A \in \mathbb{R}^{n \times n}$ and that Jacobi's method has produced an orthogonal matrix R and a symmetric matrix B such that $B = R^{-1} A R$. Suppose also that $b_{ii} < b_{jj}$ for all $i \neq j$. Show that, for each $j = 1, 2, \dots, n$, there is at least one eigenvalue of A such that

$$b_{jj} < \lambda \leq \bar{b}_j.$$

- 5.13 Suppose that $A \in \mathbb{R}^{n \times n}$ and that the Householder reduction and QR algorithm have produced an orthogonal matrix Q and a tridiagonal matrix T such that $T = Q^{-1} A Q$. Suppose also that $t_{ii} < t_{jj}$. Show that there is at least one eigenvalue of A such that

$$t_{jj} < \lambda \leq t_{jj}.$$

Polynomial interpolation

6.1 Introduction

It is time to take a break from solving equations. In this chapter we consider the problem of polynomial interpolation; it involves finding a polynomial that agrees exactly with some information that we have about a real-valued function f of a single real variable x . This information may be in the form of values $f(x_0), \dots, f(x_n)$ of the function f at some finite set of points x_0, \dots, x_n on the real line, and the corresponding polynomial is then called the **Lagrange interpolation polynomial** or, provided that f is differentiable, it may include values of the derivative of f at these points, in which case the associated polynomial is referred to as a **Hermite interpolation polynomial**.

Why should we be interested in constructing Lagrange or Hermite interpolation polynomials? If the function values $f(x)$ are known for all x in a closed interval of the real line, then the aim of polynomial

1 K \$: : #. #* B > K 6 #&*? C B
) 6D C # % #. #* 9 , D 1 , , \$
 , 15 4 O / , 1' , , \$
 %, 6 1 , , #&??4 / 6 , 6 : " \$
 6 6 1 6 6 1 , , 6 R , 5
 6 1 %, - , 4) #&. & ' 5
 : 1 O , 5 1 #. .
 * % #. #* " 1 , 5 F 8 1 H
) - 4
 2) O B / , #. / + : 9 , C # K 6 # @ #
 9 , D O 7 6 1 8
 E5 6 , 1 , 4 O , 6 G 6
 1 E1 , 4) #. &* E 1 : ,
 4 1 O :
 #. . 4 % 1
 O || 2 O 6 O || J , G
 O || 1 1 O , 4

interpolation is to approximate the function f by a polynomial over this interval. Given that any polynomial can be completely specified by its (finitely many) coefficients, storing the interpolation polynomial for f in a computer will be, generally, more economical than storing f itself.

Frequently, it is the case, though, that the function values $f(x)$ are only known at a finite set of points x_0, \dots, x_n , perhaps as the results of some measurements. The aim of polynomial interpolation is then to attempt to reconstruct the unknown function f by seeking a polynomial p whose graph in the (x, y) -plane passes through the points with coordinates $(x_i, f(x_i))$, $i = 0, \dots, n$. Of course, in general, the resulting polynomial p will differ from f (unless f itself is a polynomial of the same degree as p), so an error will be incurred. In this chapter we shall also establish results which provide bounds on the size of this error.

6.2 Lagrange interpolation

Given that n is a nonnegative integer, let \mathcal{P}_n denote the set of all (real-valued) polynomials of degree $\leq n$ defined over the set \mathcal{D} of real numbers. The simplest interpolation problem can be stated as follows: given x_0, \dots, x_n and y_0, \dots, y_n in \mathcal{D} , find a polynomial $p \in \mathcal{P}_n$ such that $p(x_i) = y_i$. The solution to this is, trivially, $p(x) = y$. The purpose of this section is to explore the following more general problem.

Let $n \geq 1$, and suppose that x_0, \dots, x_n are *distinct* real numbers (i.e., $x_i \neq x_j$ for $i \neq j$) and y_0, \dots, y_n are real numbers; we wish to find $p \in \mathcal{P}_n$ such that $p(x_i) = y_i$, $i = 0, 1, \dots, n$.

To prove that this problem has a unique solution, we begin with a useful lemma.

Lemma 6.1 *Suppose that $n \geq 1$. There exist polynomials $L_0, \dots, L_n \in \mathcal{P}_n$, $k = 0, 1, \dots, n$, such that*

$$L_k(x_i) = \begin{cases} 1, & i = k, \\ 0, & i \neq k, \end{cases} \quad (6.1)$$

for all $i, k = 0, 1, \dots, n$. Moreover,

$$p(x) = \sum_{k=0}^n L_k(x)y_k \quad (6.2)$$

satisfies the above interpolation conditions; in other words, $p \in \mathcal{P}_n$ and $p(x_i) = y_i$, $i = 0, 1, \dots, n$.

Proof For each fixed k , $0 \leq k \leq n$, L is required to have n zeros — x_i , $i = 0, 1, \dots, n, i \neq k$; thus, $L(x)$ is of the form

$$L(x) = C \prod_{i=0, i \neq k}^n (x - x_i), \quad (6.3)$$

where C

Proof In view of Remark 6.1, for $n = 0$ the proof is trivial. Let us therefore suppose that $n \geq 1$. It follows immediately from Lemma 6.1 that the polynomial p & defined by

$$p(x) = \sum_{i=0}^n L_i(x)y_i$$

satisfies the conditions (6.5), thus showing the *existence* of the required polynomial. It remains to show that p is the *unique* polynomial in \mathcal{P}_n & satisfying the interpolation property

$$p(x_i) = y_i, \quad i = 0, 1, \dots, n.$$

Suppose, otherwise, that there exists q &, different from p , such that $q(x_i) = y_i$, $i = 0, 1, \dots, n$. Then, $p - q$ & and $p - q$ has $n + 1$ distinct roots, x_i , $i = 0, 1, \dots, n$; since a polynomial of degree n cannot have more than n distinct roots, unless it is identically 0, it follows that

$$p(x) - q(x) = 0,$$

which contradicts our assumption that p and q are distinct. Hence, there exists only one polynomial p & which satisfies (6.5). \square

Definition 6.1 Suppose that $n \geq 0$. Let x_i , $i = 0, \dots, n$, be distinct real numbers, and y_i , $i = 0, \dots, n$, real numbers. The polynomial p defined by

$$p(x) = \sum_{i=0}^n L_i(x)y_i, \quad (6.6)$$

with $L_k(x)$, $k = 0, 1, \dots, n$, defined by (6.4) when $n \geq 1$, and $L_0(x) = 1$ when $n = 0$, is called the **Lagrange interpolation polynomial** of degree n for the set of points (x_i, y_i) : $i = 0, \dots, n$. The numbers x_i , $i = 0, \dots, n$, are called the **interpolation points**.

Frequently, the real numbers y_i are given as the values of a real-valued function f , defined on a closed real interval $[a, b]$, at the (distinct) interpolation points $x_i \in [a, b]$, $i = 0, \dots, n$.

Definition 6.2 Let $n \geq 0$. Given the real-valued function f , defined and continuous on a closed real interval $[a, b]$, and the (distinct) interpolation points $x_i \in [a, b]$, $i = 0, \dots, n$, the polynomial p defined by

$$p(x) = \sum_{i=0}^n L_i(x) f(x_i) \quad (6.7)$$

is the **Lagrange interpolation polynomial of degree n (with interpolation points $x_i, i = 0, \dots, n$) for the function f .**

Example 6.1 We shall construct the Lagrange interpolation polynomial of degree 2 for the function $f: x \mapsto e^{-x}$ on the interval $[-1, 1]$, with interpolation points $x_0 = -1, x_1 = 0, x_2 = 1$.

As $n = 2$, we have that

$$L_0(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} = -x(x - 1).$$

Similarly, $L_1(x) = 1 - x^2$ and $L_2(x) = -x(x + 1)$. Therefore,

$$p(x) = -x(x - 1)e^{-1} + (1 - x^2)e^0 + -x(x + 1)e^1.$$

Thus, after some simplification, $p(x) = 1 + x \sinh 1 + x^2 (\cosh 1 - 1)$.

Although the values of the function f and those of its Lagrange interpolation polynomial coincide at the interpolation points, $f(x)$ may be quite different from $p(x)$ when x is *not* an interpolation point. Thus, it is natural to ask just how large the difference $f(x) - p(x)$ is when $x = x_i, i = 0, \dots, n$. Assuming that the function f is sufficiently smooth, an estimate of the size of the **interpolation error** $f(x) - p(x)$ is given in the next theorem.

Theorem 6.2 Suppose that $n \geq 0$, and that f is a real-valued function, defined and continuous on the closed real interval $[a, b]$, such that the derivative of f of order $n + 1$ exists and is continuous on $[a, b]$. Then, given that $x \in [a, b]$, there exists $\xi = \xi(x)$ in (a, b) such that

$$f(x) - p(x) = \frac{f^{(n+1)}(\xi)}{(n+1)!} \omega(x), \quad (6.8)$$

where

$$\omega(x) = (x - x_0) \dots (x - x_n). \quad (6.9)$$

Moreover

$$|f(x) - p(x)| \leq \frac{M}{(n+1)!} |\omega(x)|, \quad (6.10)$$

where

$$M = \max_{(x, y) \in \dots} f$$

It is perhaps worth noting that since the location of x in the interval $[a, b]$ is unknown (to the extent that the exact dependence of M on x is not revealed by the proof of Theorem 6.2), (6.8) is of little practical value; on the other hand, given the function f , an upper bound on the maximum value of f over $[a, b]$ is, at least in principle, possible to obtain, and thereby we can provide an upper bound on the size of the interpolation error by means of inequality (6.10).

6.3 Convergence

An important theoretical question is whether or not a sequence (p_n) of interpolation polynomials for a continuous function f converges to f as $n \rightarrow \infty$. This question needs to be made more specific, as p_n depends on the distribution of the interpolation points $x_j, j = 0, 1, \dots, n$, not just on the value of n . Suppose, for example, that we agree to choose equally spaced points, with

$$x_j = a + j(b - a)/n, \quad j = 0, 1, \dots, n, \quad n \geq 1.$$

The question of convergence then clearly depends on the behaviour of M_n as n increases. In particular, if

$$\lim_{n \rightarrow \infty} \frac{M_n}{(n+1)!} \max_{x \in [a, b]} |f^{(n+1)}(x)| = 0,$$

then, by (6.10),

$$\lim_{n \rightarrow \infty} \max_{x \in [a, b]} |f(x) - p_n(x)| = 0, \quad (6.12)$$

and we say that the sequence of interpolation polynomials (p_n) , with equally spaced points on $[a, b]$, converges to f as $n \rightarrow \infty$, uniformly on the interval $[a, b]$.

You may now think that if all derivatives of f exist and are continuous on $[a, b]$, then (6.12) will hold. Unfortunately, this is not so, since the sequence

$$M_n = \max_{x \in [a, b]} |f^{(n+1)}(x)|$$

may tend to ∞ , as $n \rightarrow \infty$, faster than the sequence $(1/(n+1)!)$ tends to 0.

In order to convince you of the existence of such ‘pathological’ functions, we consider the sequence of Lagrange interpolation polynomials

Table 6.1. *Runge phenomenon: n denotes the degree of the interpolation polynomial p to f, with equally spaced points on [-5, 5].*

'Max error' signifies $\max_{-5 \leq x \leq 5} |f(x) - p(x)|$.

E	F!
	%(
,	"
%	%
	,
	\$
	&%%
,) (
%	' (
)'
	((\$
,	()

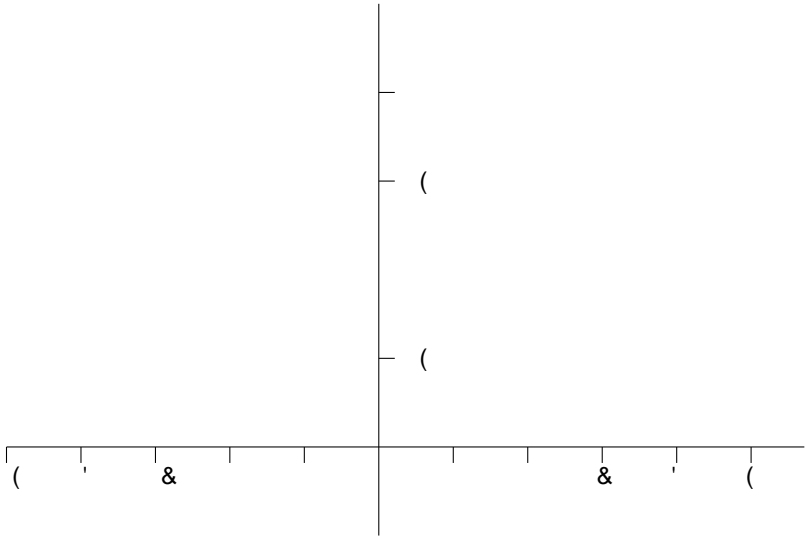
$p_n, n = 0, 1, 2, \dots$, with equally spaced interpolation points on the interval $[-5, 5]$, to

$$f(x) = \frac{1}{1+x^2}, \quad x \in [-5, 5].$$

This example is due to Runge, and the characteristic behaviour exhibited by the sequence of interpolation polynomials p_n in Table 6.1 is referred to as the **Runge phenomenon**: Table 6.1 shows the maximum difference between $f(x)$ and $p_n(x)$ for $-5 \leq x \leq 5$, for values of n from 2 up to 24. The numbers indicate clearly that the maximum error increases exponentially as n increases. Figure 6.1 shows the interpolation polynomial p_{24} , using the equally spaced interpolation points $x_j = -5 + j, j = 0, 1, \dots, 10$. The sizes of the local maxima near $x = \pm 5$ grow exponentially as the degree n increases.

Note that, in many ways, the function f is well behaved; all its deriva-

1 / 5 - B* % #.>? F 6 C * K 6 #@ &
 FI , 4 O F 6D #.. 6 , 1 J 5 6 1' \$
 F 6 , 1 , 2 5 6!
 8 , , E 1 1 , G , S, #. &
 " 5 1 , 1 , 1 4) #@
 , 5 , 2 1 1% , FI 4 O E
 4 % 1 6 J 1 , " 6 1



The construction is similar to that of the Lagrange interpolation polynomial, but now requires two sets of polynomials H and K with $k = 0, \dots, n$; these will be defined in the proof of the next theorem.

Theorem 6.3 (Hermite Interpolation Theorem) *Let $n \geq 0$, and suppose that $x_i, i = 0, \dots, n$, are distinct real numbers. Then, given two sets of real numbers $y_i, i = 0, \dots, n$, and $z_i, i = 0, \dots, n$, there is a unique polynomial p in \mathcal{P}_{2n+1} such that*

$$p(x_i) = y_i, \quad p'(x_i) = z_i, \quad i = 0, \dots, n. \tag{6.13}$$

Proof Let us begin by supposing that $n \geq 1$. As in the case of Lagrange interpolation, we start by constructing a set of auxiliary polynomials; we consider the polynomials H_k and $K_k, k = 0, 1, \dots, n$, defined by

$$\begin{aligned} H_k(x) &= [L_k(x)] (1 - 2L_k'(x)(x - x_k)), \\ K_k(x) &= [L_k(x)] (x - x_k), \end{aligned} \tag{6.14}$$

where

$$L_k(x) = \prod_{i \neq k} \frac{x - x_i}{x_k - x_i}.$$

Clearly H_k and $K_k, k = 0, 1, \dots, n$, are polynomials of degree $2n + 1$. It is easy to see that $H_k(x_i) = K_k(x_i) = 0, H_k'(x_i) = K_k'(x_i) = 0$ whenever $i, k = 0, 1, \dots, n$ and $i \neq k$; moreover, a straightforward calculation verifies their values when $i = k$, showing that

$$\begin{aligned} H_k(x_i) &= \begin{cases} 1, & i = k, \\ 0, & i \neq k, \end{cases} & H_k'(x_i) &= 0, & i, k &= 0, 1, \dots, n, \\ K_k(x_i) &= 0, & K_k'(x_i) &= \begin{cases} 1, & i = k, \\ 0, & i \neq k, \end{cases} & i, k &= 0, 1, \dots, n. \end{aligned}$$

We deduce that

$$p(x) = \sum_{k=0}^n [H_k(x)y_k + K_k(x)z_k]$$

satisfies the conditions (6.13), and p is clearly an element of \mathcal{P}_{2n+1} .

To show that this is the only polynomial in \mathcal{P}_{2n+1} satisfying these conditions, we suppose otherwise; then, there exists a polynomial q in \mathcal{P}_{2n+1} , distinct from p , such that

$$q(x_i) = y_i \quad \text{and} \quad q'(x_i) = z_i, \quad i = 0, 1, \dots, n.$$

Consequently, $p - q$ has $n + 1$ distinct zeros; therefore, Rolle's Theorem implies that, in addition to the $n + 1$ zeros $x_i, i = 0, 1, \dots, n$, $p - q$ vanishes at another n points which interlace the x_i . Hence $p - q$ & has $2n + 1$ zeros, which means that $p - q$ is identically zero, so that $p - q$ is a constant function. However, $(p - q)(x_i) = 0$ for $i = 0, 1, \dots, n$, and hence $p - q = 0$, contradicting the hypothesis that p and q are distinct. Thus, p is unique.

When $n = 0$, we define $H(x) = 1$ and $K(x) = x - x_0$, which correspond to taking $L(x) = 1$ in (6.15). Clearly, p defined by

$$p(x) = H(x)y + K(x)z = y + (x - x_0)z$$

is the unique polynomial in \mathcal{P}_1 such that $p(x_0) = y$ and $p(x_1) = z$. \square

Definition 6.3 Let $n \geq 0$, and suppose that $x_i, i = 0, \dots, n$, are distinct real numbers and $y, z, i = 0, \dots, n$, are real numbers. The polynomial p defined by

$$p(x) = [H(x)y + K(x)z] \quad (6.15)$$

where $H(x)$ and $K(x)$ are defined by (6.15), is called the **Hermite interpolation polynomial** of degree $2n + 1$ for the set of values given in $(x_i, y, z): i = 0, \dots, n$.

Example 6.2 We shall construct a cubic polynomial p such that

$$p(0) = 0, \quad p(1) = 1, \quad p'(0) = 1 \quad \text{and} \quad p'(1) = 0.$$

Here $n = 1$, and since $p(0) = p(1) = 0$ the polynomial simplifies to

$$p(x) = H(x) + K(x).$$

We easily find that, with $n = 1, x_0 = 0$ and $x_1 = 1$,

$$L(x) = 1 - x, \quad L(x) = x,$$

and then,

$$H(x) = [L(x)](1 - 2L(x)(x - x_1)) = x(3 - 2x),$$

$$K(x) = [L(x)](x - x_0) = (1 - x)x.$$

These yield the required Hermite interpolation polynomial,

$$p(x) = x + x + x.$$

Definition 6.4 Suppose that f is a real-valued function, defined on the closed interval $[a, b]$ of \mathbb{R} , and that f is continuous and differentiable on this interval. Suppose, further, that $n \geq 0$ and that $x_i, i = 0, \dots, n$, are distinct points in $[a, b]$. Then, the polynomial p defined by

$$p(x) = [H(x)f(x) + K(x)f'(x)] \tag{6.16}$$

is the **Hermite interpolation polynomial of degree $2n + 1$ with interpolation points $x_i, i = 0, \dots, n$, for f** . It satisfies the conditions

$$p(x_i) = f(x_i), \quad p'(x_i) = f'(x_i), \quad i = 0, \dots, n.$$

Pictorially, the graph of p touches the graph of the function f at the points $x_i, i = 0, \dots, n$.

To conclude this section we state a result, analogous to Theorem 6.2, concerning the error in Hermite interpolation.

Theorem 6.4 Suppose that $n \geq 0$ and let f be a real-valued function, defined, continuous and $2n + 2$ times differentiable on the interval $[a, b]$, such that $f^{(2n+2)}$ is continuous on $[a, b]$. Further, let p denote the Hermite interpolation polynomial of f defined by (6.16). Then, for each $x \in [a, b]$ there exists $\xi = \xi(x)$ in (a, b) such that

$$f(x) - p(x) = \frac{f^{(2n+2)}(\xi)}{(2n+2)!} [\omega(x)]^2, \tag{6.17}$$

where $\omega(x)$ is as defined in (6.9). Moreover,

$$|f(x) - p(x)| \leq \frac{M}{(2n+2)!} [\omega(x)]^2, \tag{6.18}$$

where $M = \max_{\xi \in [a, b]} |f^{(2n+2)}(\xi)|$.

Proof The inequality (6.18) is a straightforward consequence of (6.17). In order to prove (6.17), we observe that it is trivially true if $x = x_i$

for some i , $i = 0, \dots, n$; thus, it suffices to consider $x \in [a, b]$ such that $x = x_i$, $i = 0, \dots, n$. For such x , let us define the function

$x_i, i = 0, 1, \dots, n$, are distinct points in $[a, b]$, and that $p_n(x)$ is the Lagrange interpolation polynomial for f defined by these points. Then, there exist distinct points $\xi_i, i = 1, \dots, n$, in (a, b) , and corresponding to each x in $[a, b]$ there exists a point $\xi = \xi(x)$ in (a, b) , such that

$$f(x) - p_n(x) = \frac{f^{(n)}(\xi)}{n!} \omega(x), \tag{6.20}$$

where

$$\omega(x) = (x - x_1) \dots (x - x_n).$$

Proof Since $f(x_i) - p_n(x_i) = 0, i = 0, 1, \dots, n$, there exists a point in (x_{i-1}, x_i) at which $f^{(i)}(\xi_i) - p_n^{(i)}(\xi_i) = 0$, for each $i = 1, \dots, n$. This defines the points $\xi_i, i = 1, \dots, n$. Now the proof closely follows that of Theorem 6.2.

When $x = x_i$ for some $i = 1, \dots, n$, both sides of (6.20) are zero. Suppose then that x is distinct from all the $x_i, i = 1, \dots, n$, and define the function $t = t(x)$ by

$$t = t(x) = p_n(t) - \frac{f(x) - p_n(x)}{\omega(x)} \omega(t).$$

This function vanishes at every point $x_i, i = 1, \dots, n$, and also at the point $t = x$. By successively applying Rolle's Theorem we deduce that $t^{(i)} = 0$ vanishes at some point ξ_i . The result then follows as in the proof of Theorem 6.2. □

Corollary 6.1 *Under the conditions of Theorem 6.5,*

$$|f(x) - p_n(x)| \leq \frac{M}{n!} \omega(x) = \frac{(b-a)^{n+1}}{n!} M$$

for all x in $[a, b]$, where $M = \max_{x \in [a, b]} |f^{(n+1)}(x)|$.

In particular, we deduce that if f and all its derivatives are defined and continuous on the closed interval $[a, b]$, and

$$\lim_{n \rightarrow \infty} \frac{(b-a)^{n+1}}{n!} M = 0,$$

then $\lim_{n \rightarrow \infty} \max_{x \in [a, b]} |f(x) - p_n(x)| = 0$, showing the convergence of the sequence of interpolation polynomials (p_n) to f , uniformly on $[a, b]$.

The discussion in the last few paragraphs may give the impression that numerical differentiation is a straightforward procedure. In practice, however, things are much more complicated since the function values $f(x_i), i = 0, 1, \dots, n$, will be polluted by rounding errors.

Example 6.3

and $x = h$, tends to 0; in the latter case, the degree of the polynomial p tends to infinity and consequently the spacing between the increasing number of consecutive interpolation points shrinks. Nevertheless, Example 6.3 illustrates the issue that caution should be exercised in the course of numerical differentiation when rounding errors are present.

6.6 Notes

The interpolation polynomial (6.6) was discovered by Edward Waring (1736–1798) in 1776, rediscovered by Euler in 1783 and published by Joseph-Louis Lagrange (1736–1813) in his *Leçons élémentaires sur les mathématiques*, Paris, 1795.

Lagrange's interpolation theorem is a purely algebraic result, and it also holds in number fields different from the field of real numbers considered in this chapter. In particular, it holds if the numbers x and y , $i = 0, 1, \dots, n$, are complex, and the polynomial p has complex coefficients. Theorem 6.2 is due to Augustin-Louis Cauchy (1789–1857). The interpolation polynomial (6.15) was discovered by Charles Hermite (1822–1901).

Before modern computers came into general use about 1960, the evaluation of a standard mathematical function for a given value of x required the use of published tables of the function, in book form. If x was not one of the tabulated values, the required result was obtained by interpolation, using tabulated values close to x . The tabulated values were given at equally spaced points, so that usually $x = jh$, where h is a fixed increment. In this case the Lagrange formula can be simplified; as this sort of interpolation had to be done frequently, various devices were used to make the calculations easy and quick. Older books, such as F.B. Hildebrand's *Introduction to Numerical Analysis*, published in 1956, contain extensive discussions of such special methods of interpolation, some of which date back to the time of Newton, but are now mainly of historical interest. A notable early contribution to the development of mathematical tables is the work of Henry Briggs (1560–1630), Savilian Professor of Geometry and fellow of Merton College in Oxford, entitled *Arithmetica logarithmica*, published in 1624. It contained extensive calculations of the logarithms of thirty thousand numbers to 14 decimal digits; these were the numbers from 1 to 20000 and from 90000 to 100000. It also contained tables of the sin function to 15 decimal digits, and of the tan and sec functions to 10 decimal digits.

Exercises

- 6.1 Construct the Lagrange interpolation polynomial p of degree 1, for a continuous function f defined on the interval $[-1, 1]$, using the interpolation points $x = -1, x = 1$. Show further that if the second derivative of f exists and is continuous on $[0, 1]$, then

$$f(x) - p(x) = \frac{M}{2}(1 - x^2) = \frac{M}{2}, \quad x \in [-1, 1],$$

where $M = \max_{x \in [0, 1]} |f''(x)|$. Give an example of a function f , and a point x , for which equality is achieved.

- 6.2 (i) Write down the Lagrange interpolation polynomial of degree 1 for the function $f: x \rightarrow x^2$, using the points $x = 0, x = a$. Verify Theorem 6.2 by direct calculation, showing that in this case p is unique and has the value $p(x) = -x(x + a)$.
 (ii) Repeat the calculation for the function $f: x \rightarrow (2x - a)^2$; show that in this case there are two possible values for p , and give their values.

- 6.3 Given the distinct points $x_i, i = 0, 1, \dots, n + 1$, and the points $y_i, i = 0, 1, \dots, n + 1$, let q be the Lagrange polynomial of degree n for the set of points $(x_i, y_i): i = 0, 1, \dots, n$ and let r be the Lagrange polynomial of degree n for the points $(x_i, y_i): i = 1, 2, \dots, n + 1$. Define

$$p(x) = \frac{(x - x_{n+1})r(x) - (x - x_0)q(x)}{x - x_0 - x_{n+1}}.$$

Show that p is the Lagrange polynomial of degree $n + 1$ for the points $(x_i, y_i): i = 0, 1, \dots, n + 1$.

- 6.4 Let $n \geq 1$. The points x_j are equally spaced in $[-1, 1]$, so that

$$x_j = \frac{2j - n}{n}, \quad j = 0, \dots, n.$$

With the usual notation

$$\omega(x) = (x - x_0) \dots (x - x_n),$$

show that

$$(1 - 1/n) \omega(x) = \frac{(2n)!}{2^n n! n!}.$$

Using Stirling's formula

$$N! \sim \sqrt{2\pi N} e^{-N} N^N,$$

verify that

$$(1 - 1/n)^{2n} \sim \frac{2}{n} e^{-2}$$

for large values of n .

6.5 Let $n \geq 1$. Suppose that $x_i, i = 0, 1, \dots, n$, are distinct real numbers, and $y_i, u_i, i = 0, 1, \dots, n$, are real numbers. Suppose, further, that there exists $p_1(x)$ & $p_2(x)$ such that $p_1(x_i) = y_i$ for all $i = 0, 1, \dots, n$, and $p_2(x_i) = u_i, i = 0, 1, \dots, n$. Attempt to prove that $p_1(x)$ is the unique polynomial with these properties, by adapting the uniqueness proofs in Sections 6.2 and 6.4, using Rolle's Theorem; explain where the proof fails. Show that there is no polynomial $p(x)$ & $q(x)$ such that $p(x_i) = 1, p(x_0) = 0, p(x_1) = 1, p(x_{n-1}) = 0, p(x_n) = 0, p(x_1) = 0$, but that if the first condition is replaced by $p(x_{n-1}) = 1$, then there is an infinite number of such polynomials. Give an explicit expression for the general form of these polynomials.

6.6 Suppose that $n \geq 1$. The function f and its derivatives of order up to and including $2n + 1$ are continuous on $[a, b]$. The points $x_i, i = 0, 1, \dots, n$, are distinct and lie in $[a, b]$. Construct polynomials $l(x), h(x), k(x), i = 1, \dots, n$, of degree $2n$ such that the polynomial

$$p(x) = l(x)f(x) + \sum_{i=1}^n [h_i(x)f(x) + k_i(x)f'(x)]$$

satisfies the conditions

$$p(x_i) = f(x_i), \quad i = 0, 1, \dots, n,$$

and

$$p'(x_i) = f'(x_i), \quad i = 1, \dots, n.$$

Show also that for each value of x in $[a, b]$ there is a number δ , depending on x , such that

$$f(x) - p(x) = \frac{(x - x_0)^{2n+1}}{(2n+1)!} f^{(2n+1)}(\xi).$$

6.7 Suppose that $n \geq 2$. The function f and its derivatives of order up to and including $2n$ are continuous on $[a, b]$. The points $x_i, i = 0, 1, \dots, n$, are distinct and lie in $[a, b]$. Explain how to

construct polynomials $l_1(x), l_2(x), h(x), k(x), i = 1, \dots, n-1$, of degree $2n-1$ such that the polynomial

$$p_{2n-1}(x) = l_1(x)f_1(x) + l_2(x)f_2(x) + \dots + [h(x)f(x) + k$$

- Find the limit of this expression as $h \rightarrow 0$, and deduce that $p_n(x) - q_n(x) \rightarrow 0$ as $h \rightarrow 0$, where $q_n(x)$ is the Hermite interpolation polynomial for f , using the points $x_i, i = 0, \dots, n-1$.
- 6.10 Construct the Hermite interpolation polynomial of degree 3 for the function $f: x \rightarrow x^2$, using the points $x_0 = 0, x_1 = a$, and show that it has the form $p(x) = 3a x^2 - 2a x$. Verify Theorem 6.4 by direct calculation, showing that in this case p is unique and has the value $p'(x) = -(x + 2a)$.
- 6.11 The complex function $z \rightarrow f(z)$ of the complex variable z is holomorphic in the region D of the complex plane; the boundary of D is the simple closed contour C . The interpolation points $x_j, j = 0, 1, \dots, n$, with $n \geq 1$, and the point x all lie in D . Determine the residues of the function g defined by

$$g(z) = \frac{f(z)}{z-x} + \sum_{j=0}^n \frac{f(x_j)}{z-x_j} - \frac{p(x)}{z-x}$$

at its poles in D , and deduce that

$$f(x) - p(x) = \frac{1}{2\pi i} \int_C \frac{f(z)}{z-x} - \frac{p(x)}{z-x} dz,$$

where p is the Lagrange interpolation polynomial for the function f using the interpolation points $x_j, j = 0, 1, \dots, n$.

Now, suppose that the real number x and the interpolation points $x_j, j = 0, 1, \dots, n$, all lie in the real interval $[a, b]$, and that D consists of all the points z such that $z = t + iK$ for all $t \in [a, b]$, where K is a constant with $K > b - a$. Show that the length of the contour C is $2(b - a) + 2K$, and that

$$f(x) - p(x) < \frac{(b - a + K)M}{K},$$

where M is such that $|f(z)| \leq M$ on C . Deduce that the sequence (p_n) converges to f , uniformly on $[a, b]$.

- Show that these conditions are not satisfied by the function $f: x \rightarrow 1/(1+x)$ for x in the interval $[-5, 5]$. For what values of a are the conditions satisfied by f for x in the interval $[a, a]$?
- 6.12 With the same notation as in Example 6.3, let

$$E(h) = \frac{(f(h) + f(-h)) - 2f(0)}{2h}$$

Suppose that $f'(x)$ exists and is continuous at all $x \in [-h, h]$.

By expanding $f(h)$ and $f(-h)$ into Taylor series about the point 0, show that there exists $\delta(h, h)$ such that

$$E(h) = \frac{1}{6}h^3 f'''(\xi) + \frac{1}{24}h^5 f^{(5)}(\eta).$$

Hence deduce that

$$|E(h)| \leq \frac{1}{6}h^3 M + \frac{1}{24}h^5 N$$

where $M = \max_{x \in [a, b]} |f'''(x)|$ and $N = \max_{x \in [a, b]} |f^{(5)}(x)|$. Show further that the right-hand side of the last inequality achieves its minimum value when

$$h = \frac{3}{M}.$$

Numerical integration – I

7.1 Introduction

The problem of evaluating definite integrals arises both in mathematics and beyond, in many areas of science and engineering. At some point in our mathematical education we all learned to calculate simple integrals such as

$$\int_0^1 e^{-x} dx \quad \text{or} \quad \int_0^{\pi/2} \cos x dx$$

using a table of integrals, so you will know that the values of these are $e^{-1} + 1$ and 0 respectively; but how about the innocent-looking

$$\int_0^1 e^{-x^2} dx \quad \text{and} \quad \int_0^1 \cos(x^2) dx,$$

or the more exotic

$$\int_0^1 \exp(\sin(\cos(\sinh(\cosh(\tan^{-1}(\log(x))))) dx?$$

Please try to evaluate these using a table of integrals and see how far you can get! It is not so simple, is it? Of course, you could argue that the last example was completely artificial. Still, it illustrates the point that it is relatively easy to think of a continuous real-valued function f defined on a closed interval $[a, b]$ of the real line such that the definite integral

$$\int_a^b f(x) dx \tag{7.1}$$

points; for the sake of simplicity, we shall assume that these are equally spaced, that is,

$$x = a + ih, \quad i = 0, 1, \dots, n,$$

where

$$h = (b - a)/n.$$

The Lagrange interpolation polynomial of degree n for the function f , with these interpolation points, is of the form

$$p(x) = \sum_{k=0}^n L_k(x)f(x_k) \quad \text{where} \quad L_k(x) = \prod_{j \neq k} \frac{x - x_j}{x_k - x_j}.$$

Inserting the expression for p into the right-hand side of (7.2) yields

$$\int_a^b f(x) dx \approx \sum_{k=0}^n w_k f(x_k), \quad (7.3)$$

where

$$w_k = \int_a^b L_k(x) dx, \quad k = 0, 1, \dots, n. \quad (7.4)$$

The values w_k , $k = 0, 1, \dots, n$, are referred to as the **quadrature weights**, while the interpolation points x_k , $k = 0, 1, \dots, n$, are called the **quadrature points**. The numerical quadrature rule (7.3), with quadrature weights (7.4) and equally spaced quadrature points, is called the **Newton–Cotes formula** of order n . In order to illustrate the general idea, we consider two simple examples.

Trapezium rule. In this case we take $n = 1$, so that $x_0 = a$, $x_1 = b$; the Lagrange interpolation polynomial of degree 1 for the function f is simply

$$\begin{aligned} p(x) &= L_0(x)f(a) + L_1(x)f(b) \\ &= \frac{x-b}{a-b}f(a) + \frac{x-a}{b-a}f(b) \\ &= \frac{1}{b-a}[(b-x)f(a) + (x-a)f(b)]. \end{aligned}$$

Integrating $p(x)$ from a to b yields

$$\int_a^b f(x) dx \approx \frac{b-a}{2}[f(a) + f(b)].$$

This numerical integration formula is called the trapezium rule. The

terminology stems from the fact that the expression on the right is the area of the trapezium with vertices $(a, 0)$, $(b, 0)$, $(a, f(a))$, $(b, f(b))$.

Simpson's rule. A slightly more sophisticated quadrature rule is obtained by taking $n = 2$. In this case $x_0 = a$, $x_1 = (a + b)/2$ and $x_2 = b$, and the function f is approximated by a quadratic Lagrange interpolation polynomial.

The quadrature weights are calculated from

$$\begin{aligned}
 w_0 &= \int_a^b L_0(x) dx \\
 &= \int_a^b \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} dx \\
 &= \int_a^b \frac{t(t - 1)}{2} \frac{b - a}{2} dt \\
 &= \frac{b - a}{6},
 \end{aligned}$$

where it is convenient to make the change of variable

$$x = \frac{b - a}{2} t + \frac{b + a}{2}.$$

Similarly, $w_2 = (b - a)/6$, and it is easy to see that $w_0 = w_2$ by symmetry. This gives

$$\int_a^b f(x) dx \approx \frac{b - a}{6} [f(a) + 4f\left(\frac{a + b}{2}\right) + f(b)],$$

a numerical integration formula known as Simpson's rule.

It is very important to notice that the weights w_k defined in (7.4) depend only on n and k , not on the function f . Their values can therefore

n	k	w_k
1	0	1
2	0	1/6
2	1	4/6
2	2	1/6
3	0	3/8
3	1	6/8
3	2	3/8
4	0	7/90
4	1	32/90
4	2	48/90
4	3	16/90
4	4	7/90
5	0	8/252
5	1	32/252
5	2	48/252
5	3	32/252
5	4	8/252
5	5	8/252

be calculated in advance, as in the trapezium rule and Simpson’s rule. The evaluation of the approximation to the integral (7.1) is then a trivial matter; it is only necessary to compute $f(x_k)$ at each of the quadrature points x_k , $k = 0, 1, \dots, n$, multiply by the known weights w_k for $k = 0, 1, \dots, n$, and form the sum on the right-hand side of (7.3).

7.3 Error estimates

Our next task is to estimate the size of the error in the numerical integration formula (7.3), that is, the error that has been committed by integrating the interpolating Lagrange polynomial of f instead of f itself. The error in (7.3) is defined by

$$E(f) = \int_a^b f(x) dx - \sum_{k=0}^n w_k f(x_k).$$

The next theorem provides a useful bound on $E(f)$ under the additional hypothesis that the function f is sufficiently smooth.

Theorem 7.1 *Let $n \geq 1$. Suppose that f is a real-valued function, defined and continuous on the interval $[a, b]$, and let $f^{(n+1)}$ be defined and continuous on $[a, b]$. Then,*

$$|E(f)| \leq \frac{M}{(n+1)!} \int_a^b \omega(x) dx, \tag{7.5}$$

where $M = \max_{x \in [a, b]} |f^{(n+1)}(x)|$ and $\omega(x) = (x - x_0) \dots (x - x_n)$.

Proof Recalling the definition of the weights w_k from (7.4), we can write $E(f)$ as follows:

$$\begin{aligned} E(f) &= \int_a^b f(x) dx - \sum_{k=0}^n w_k f(x_k) \\ &= \int_a^b [f(x) - p_n(x)] dx. \end{aligned}$$

Thus,

$$|E(f)| \leq \int_a^b |f(x) - p_n(x)| dx.$$

The desired error estimate (7.5) follows by inserting (6.8) into the right-hand side of this inequality. □

Let us use this theorem to estimate the size of the error which arises from applying the trapezium rule to the integral $\int_a^b f(x) dx$. In this case, with $n = 1$ and $w(x) = (x-a)(x-b)$, the bound (7.5) reduces to

$$\begin{aligned} E(f) &= \frac{M}{2} \int_a^b (x-a)(x-b) dx \\ &= \frac{M}{2} \int_a^b (b-x)(x-a) dx \\ &= \frac{(b-a)^2}{12} M. \end{aligned} \quad (7.6)$$

An analogous but slightly more tedious calculation shows that, for Simpson's rule,

$$\begin{aligned} E(f) &= \frac{M}{6} \int_a^b (x-a)(x-(a+b)/2)(x-b) dx \\ &= \frac{(b-a)^3}{196} M. \end{aligned} \quad (7.7)$$

Unfortunately, (7.7) gives a considerable overestimate of the error in Simpson's rule; in particular it does not bring out the fact that $E(f) = 0$ whenever f is a polynomial of degree 3. The next theorem will allow us to give a sharper bound on the error in Simpson's rule which illustrates this fact. More generally, it is quite easy to prove that when n is odd the Newton-Cotes formula (7.3) (with w defined by (7.4)) is exact for all polynomials of degree n , while when n is even it is also exact for all polynomials of degree $n+1$ (see Exercise 2 at the end of the chapter).

Theorem 7.2 *Suppose that f is a real-valued function, defined and continuous on the interval $[a, b]$, and that $f^{(4)}$ is continuous on $[a, b]$. Then,*

$$\int_a^b f(x) dx - \frac{b-a}{6} [f(a) + 4f((a+b)/2) + f(b)] = \frac{(b-a)^5}{2880} f^{(4)}(\xi), \quad (7.8)$$

for some ξ in (a, b) .

Proof Making the change of variable

$$x = \frac{a+b}{2} + \frac{b-a}{2}t, \quad t \in [-1, 1],$$

and defining the function $t \rightarrow F(t)$ by $F(t) = f(x)$, we see that

$$\begin{aligned} \int_a^b f(x) dx &= \frac{b-a}{6} [f(a) + 4f((a+b)/2) + f(b)] \\ &= \frac{b-a}{2} \int_{-1}^1 F(t) dt = \frac{1}{3} [F(-1) + 4F(0) + F(1)]. \end{aligned} \quad (7.9)$$

We now introduce the function $t \rightarrow G(t)$ by

$$G(t) = \int_{-1}^t F(s) ds = \frac{t}{3} [F(-t) + 4F(0) + F(t)], \quad t \in [-1, 1];$$

the right-hand side of (7.9) is then simply $-(b-a)G(1)$.

The remainder of the proof is devoted to showing that $-(b-a)G(1)$ is, in turn, equal to the right-hand side of (7.8) for some ξ in (a, b) . To do so, we define

$$H(t) = G(t) - tG(1), \quad t \in [-1, 1],$$

and apply Rolle's Theorem repeatedly to the function H . Noting that $H(0) = H(1) = 0$, we deduce that there exists $\xi_1 \in (0, 1)$ such that $H'(\xi_1) = 0$. But it is easy to show that $H'(0) = 0$, so there exists $\xi_2 \in (0, \xi_1)$ such that $H'(\xi_2) = 0$. Again we see that $H'(0) = 0$, so there exists $\xi_3 \in (0, \xi_2)$ such that $H'(\xi_3) = 0$. Now,

$$G'(t) = \frac{t}{3} [F'(t) - F'(-t)],$$

and therefore

$$H'(\xi_3) = \frac{1}{3} [F'(\xi_3) - F'(-\xi_3)] = 60 G(1).$$

Applying the Mean Value Theorem to the function F' this shows that there exists $\xi \in (-\xi_3, \xi_3)$ such that

$$\begin{aligned} H'(\xi) &= \frac{1}{3} [2F'(\xi)] = 60 G(1) \\ &= \frac{2}{3} [F'(\xi) + 90G(1)]. \end{aligned}$$

Since $H'(\xi) = 0$ and $G(1) \neq 0$, this means that

$$G(1) = \frac{1}{90} F'(\xi) = \frac{(b-a)}{1440} f'(\xi),$$

and the required result follows. \square

Table 7.1. 1 is the result of the Newton–Cotes formula of degree n for the approximation of the integral (7.12)

	& '%
	%)\$')
&	'(
'	&)'
(&)%\$
%	&)' (
)	\$ \$\$
	('\$
\$	&\$ %
	'%)&&
	& ''))
	& \$'
&	\$ \$
') \$\$\$('
(' (((%

7.4 The Runge phenomenon revisited

By looking at the right-hand side of the error bound (7.5) we may be led to believe that by increasing n, that is by approximating the integrand by Lagrange interpolation polynomials of increasing degree and integrating these exactly, we shall reduce the size of the quadrature error E (f). However, this is not always the case, even for very smooth functions f. An example of this behaviour uses the same function as in Section 6.3; Table 7.1 gives the results of applying Newton–Cotes formulae of increasing degree to the evaluation of the integral

$$\int_{-1}^1 \frac{1}{1+x} dx. \tag{7.12}$$

These results do not evidently converge as n increases, and in fact they eventually increase without bound. This behaviour is related to the fact that the weights w in the Newton–Cotes formula are not all positive when n > 8. We shall return to this point in Theorem 10.2.

A better approach to improving accuracy is to divide the interval [a, b] into an increasing number of subintervals of decreasing size, and then to use a numerical integration formula of fixed order n on each

of the subintervals. Quadrature rules based on this approach are called composite formulae; in the next section we shall describe two examples.

7.5 Composite formulae

We shall consider only some very simple composite quadrature rules: the composite trapezium rule and the composite Simpson rule.

Suppose that f is a function, defined and continuous on a nonempty closed interval $[a, b]$ of the real line. In order to construct an approximation to

$$\int_a^b f(x) dx,$$

we now select an integer $m \geq 2$ and divide the interval $[a, b]$ into m equal subintervals, each of width $h = (b - a)/m$, so that

$$\int_a^b f(x) dx = \sum_{i=1}^m \int_{x_{i-1}}^{x_i} f(x) dx, \tag{7.13}$$

where

$$x_i = a + ih = a + \frac{i}{m}(b - a), \quad i = 0, 1, \dots, m.$$

Each of the integrals is then evaluated by the trapezium rule,

$$\int_{x_{i-1}}^{x_i} f(x) dx \approx \frac{1}{2}h[f(x_{i-1}) + f(x_i)]; \tag{7.14}$$

summing these over $i = 1, 2, \dots, m$ leads to the following definition.

Definition 7.1 (Composite trapezium rule)

$$\int_a^b f(x) dx \approx h \left[\frac{1}{2}f(x_0) + f(x_1) + \dots + f(x_{m-1}) + \frac{1}{2}f(x_m) \right]. \tag{7.15}$$

1 , 1 , 1 6 , , " 1 \$
 1 , 5 1 \$, 7 ,
 6 5 6 " 1% , 1 6 , B . &
 6 " B) 60 C # 6 , B) 6004 % , N \$
 , 6 " 1 ! " #@@. , 1 1 , %
 F 6 ' 6 %2

The error in the composite trapezium rule can be estimated by using the error bound (7.6) for the trapezium rule on each individual subinterval $[x_{i-1}, x_i]$, $i = 1, 2, \dots, m$. For this purpose, let us define

$$\begin{aligned} \mathcal{E}_T(f) &= \int_a^b f(x) dx - h \left[\frac{1}{2}f(x_0) + f(x_1) + \dots + f(x_{m-1}) + \frac{1}{2}f(x_m) \right] \\ &= \sum_{i=1}^m \int_{x_{i-1}}^{x_i} f(x) dx - h \left[\frac{1}{2}f(x_{i-1}) + f(x_i) \right]. \end{aligned}$$

Applying (7.6) to each of the terms under the summation sign we obtain

$$\begin{aligned} \mathcal{E}_T(f) &\leq \frac{1}{12} h^3 \sum_{i=1}^m \max_{x \in [x_{i-1}, x_i]} |f''(x)| \\ &= \frac{(b-a)^3}{12m^2} M, \end{aligned} \tag{7.16}$$

where $M = \max_{x \in [a, b]} |f''(x)|$.

For Simpson’s rule, let us suppose that the interval $[a, b]$ has been divided into $2m$ intervals by the points $x_i = a + ih$, $i = 0, 1, \dots, 2m$, with $m \geq 2$ and

$$h = \frac{b-a}{2m},$$

and let us apply Simpson’s rule on each of the intervals $[x_{i-2}, x_i]$, $i = 1, 2, \dots, m$, giving

$$\begin{aligned} \int_a^b f(x) dx &= \sum_{i=1}^m \int_{x_{i-2}}^{x_i} f(x) dx \\ &= \sum_{i=1}^m \frac{2h}{6} [f(x_{i-2}) + 4f(x_{i-1}) + f(x_i)]. \end{aligned}$$

This leads to the following definition.

Definition 7.2 (Composite Simpson rule)

$$\int_a^b f(x) dx \approx \frac{h}{3} [f(x_0) + 4f(x_1) + 2f(x_2) + 4f(x_3) + \dots + 2f(x_{2m-2}) + 4f(x_{2m-1}) + f(x_{2m})]. \tag{7.17}$$

A schematic view of the pattern in which the coefficients 1, 4 and 2 appear in the composite Simpson rule is shown in Figure 7.1.



$$G = \int_a^b f(x) dx \approx \frac{h}{3} [f(x_0) + 4f(x_1) + f(x_2) + \dots + 4f(x_{m-1}) + f(x_m)]$$

In order to estimate the error in the composite Simpson rule, we proceed in the same way as for the composite trapezium rule. Let us define

$$S_j = \int_{x_{j-1}}^{x_j} f(x) dx \approx \frac{h}{3} [f(x_{j-1}) + 4f(x_{j-1/2}) + f(x_j)]$$

Applying (7.10) to each individual term in the sum and recalling that $b - a = 2mh$ we obtain the following error bound:

$$|E_S| \leq \frac{(b-a)^3}{2880m^3} M, \tag{7.18}$$

where $M = \max_{x \in [a,b]} |f'''(x)|$.

The composite rules (7.15) and (7.17) provide greater accuracy than the basic formulae considered in Section 7.2; this is clearly seen by comparing the error bounds (7.16) and (7.18) for the two composite rules with (7.6) and (7.8), the error estimates for the basic trapezium rule and Simpson rule respectively. The inequalities (7.16) and (7.18) indicate that, as long as the function f is sufficiently smooth, the errors in the composite rules can be made arbitrarily small by choosing a sufficiently large number of subintervals.

7.6 The Euler–Maclaurin expansion

We have seen in (7.16) that the error in the composite trapezium rule is bounded by a term involving $1/m^2$, where m is the number of subdivi-

sions of the interval $[a, b]$; the **Euler –Maclaurin** expansion expresses this error as a series in powers of $1/m$, and makes it possible to improve accuracy by extrapolation methods.

We first define a sequence of polynomials.

Definition 7.3 Consider the sequence of polynomials q_r , $r = 1, 2, \dots$, defined by their properties, as follows:

- (i) q_r is a polynomial of degree r ;
- (ii) for each positive integer r , $q_{r+1} = q_r$;
- (iii) q_r is an odd function if r is odd, and an even function if r is even;
- (iv) if $r > 1$ is odd, then $q_r(-1) = 0$ and $q_r(1) = 0$;
- (v) $q_0(t) = t$.

Using these conditions it is easy to construct the polynomials q_r in succession. From (v) and (ii) we get

$$q_0(t) = -t + A, \quad q_1(t) = -t + A t + A,$$

where A_0 and A_1 are constants. From (iii) we see that $A_0 = 0$; then, from (iv) it follows that $A_1 = -$. Hence,

$$q_0(t) = -t + -, \quad q_1(t) = -t + -t.$$

We can then go on to construct q_2 and q_3 , and so on.

1 : B#> % #& & + C #. #&.* \$
 1 1 D E, 1 1 1 6
 1 1 5 1 14 1 1 #*, 6 2
 , 5 , 5 6 6
 , , 15 J 6 6 J G 6
 , , 5 6 1 6 6 \$ 6 , 6
 , , 9 #&?> J 1 6 ,
 1 R 8 , 6 4 J , , N 1
 1 , 1 1 1 1 1 1 G 1 # U B D 1 U²
 1 E 1 E J , 4 C# K #& ?
 2 , , B9 6 #?@. % 6 5 6 1 F
 1 ## , 1 6 , 6 O 7 # 4) #&#@ 1 #
 , 9 1 O % " 1&?* 5
 4 , , O % " E 6 , 8 1 N
 , 1 4 " , 8

where $t = t(x) = \frac{1}{h}(x - x_{i-1})$ for $x \in [x_{i-1}, x_i]$, $i = 1, \dots, m$, and $c_r = q_r(1)/2$ for $r = 1, \dots, k$.

Proof We express the integral as a sum over the m subintervals $[x_{i-1}, x_i]$, $i = 1, \dots, m$, as in (7.13). In the interval $[x_{i-1}, x_i]$ we change the variable by writing $x = x_{i-1} + h(t+1)/2$, so that

$$f(x)dx = \frac{h}{2} g(t)dt,$$

where $f(x) = g(t)$. According to Theorem 7.3, then,

$$\begin{aligned} \int_{x_{i-1}}^{x_i} f(x)dx &= \frac{h}{2} [f(x_{i-1}) + f(x_i)] \\ &= \frac{h}{2} \int_{-1}^1 g(t)dt = \frac{h}{2} [g(-1) + g(1)] \\ &= \frac{h}{2} \int_{-1}^1 q(t)g(t)dt. \end{aligned}$$

On noting that $g(t) = (h/2) f(x)$, $i = 1, 2, \dots, 2k$, $dt = (2/h) dx$, summation over all the subintervals $[x_{i-1}, x_i]$, for $i = 1, \dots, m$, gives the required result. The important point is the symmetry of the polynomials q_r , which ensures that $q_r(1) = q_r(-1)$, so that all the derivatives of f at the internal points x cancel in the course of summation, leaving only the derivatives at a and b . □

Remark 7.1 *By successively computing the polynomials $q_r(t)$, we can determine the values of $c_r = q_r(1)/2$, $r = 1, 2, 3, \dots$. For example,*

$$c_1 = \frac{1}{4}, c_2 = \frac{1}{8}, c_3 = \frac{1}{16}, c_4 = \frac{1}{32}, c_5 = \frac{1}{64}, \dots$$

It can be shown that $c_r = \frac{B_{2r}}{2^{2r}}$ for all $r = 1, 2, 3, \dots$, where B_{2r} are the Bernoulli numbers with even index, which can be determined from

$$\sum_{k=0}^{2r-1} \binom{2r-1}{k} B_k = 0, \quad B_0 = 1, \quad B_1 = -\frac{1}{2}, \quad B_2 = \frac{1}{6}, \quad B_4 = -\frac{1}{30}, \quad B_6 = \frac{1}{42}, \dots$$

the Taylor series expansion

$$-\coth^{-1} \frac{x}{a} = \frac{B(x)}{(2r)!}.$$

Easier still, typing `integrate(f(x), x=a, x=b, method='taylor', order=2r)` at the Maple command line gives $C = \frac{f(b) - f(a)}{2r!} \dots$; c_0, c_1, \dots can be found in the same way.

An interesting consequence of Theorem 7.4 concerns the numerical integration of smooth periodic functions. Suppose that f is a continuous function defined on (a, b) such that all derivatives of f , up to and including order $2k$, are defined and continuous on (a, b) , and f is periodic on (a, b) with period $b - a$; i.e., $f(x + b - a) = f(x) = 0$ for all $x \in (a, b)$. Hence, by successive differentiation of this equality and taking $x = a$ we deduce that, in particular,

$$f^{(2r)}(b) - f^{(2r)}(a) = 0 \quad \text{for } r = 1, 2, \dots, k.$$

Therefore, according to (7.20), we have that

$$|T(m)| = O(h^{2k}).$$

The fact that for $k \geq 1$ this integration error is much smaller than the $O(h)$ error that will be observed in the case of a nonperiodic function indicates that the composite trapezium rule is particularly well suited for the numerical integration of smooth periodic functions.

A second application of the Euler–Maclaurin expansion concerns extrapolation methods. This subject will be discussed in the next section.

7.7 Extrapolation methods

In general the calculation of the higher derivatives involved in the Euler–Maclaurin expansion (7.20) is not possible. However, the existence of the expansion allows us to eliminate successive terms by repeated calculation of the trapezium rule approximation.

For example, the case $k = 2$ of (7.20) may be written in the form

$$\int_a^b f(x) dx - T(m) = C h^3 + O(h^5),$$

where $C = \frac{1}{6} [f''(b) - f''(a)]$ and $h = (b - a)/m$. This also means that

$$\int_a^b f(x) dx - T(2m) = C (h/2)^3 + O(h^5).$$

We can eliminate the term in h from these two equalities, giving

$$\int f(x)dx = \frac{4T(2m) - T(m)}{3} + O(h^3).$$

The same elimination process could be used for any two values of m , from the calculation of $T(m)$ and $T(2m)$; the advantage of using m and $2m$ is that in the computation of $T(2m)$ half the required values of $f(x)$ are already known from $T(m)$, and we do not have to calculate them again. This process of eliminating the term in h from the expansion of the error is known as **Richardson extrapolation** or **h extrapolation**. It is easy to extend the process to higher-order terms. For example,

$$\int f(x)dx = T(m) = C_1 h + C_2 h^2 + C_3 h^3 + O(h^4).$$

Hence

$$\int f(x)dx = \frac{4T(2m) - T(m)}{3} = -C_1 h + C_2 h^2 + O(h^3),$$

which leads to

$$\int f(x)dx = \frac{16T(2m) - T(m)}{15} + O(h^4),$$

where

$$T(2m) = \frac{4T(m) - T(m)}{3}.$$

Therefore,

$$T(2m) = \frac{16T(2m) - T(m)}{15}$$

approximates the integral $\int f(x)dx$ to accuracy $O(h^4)$. Adopting the notational convention

$$T(2m) = T(m)$$

and proceeding recursively,

$$T_{2^k}(2m) = \frac{2^{2k} T_{2^{k-1}}(2m) - T_{2^{k-1}}(m)}{2^{2k} - 1}$$

Table 7.2. Romberg table.

	T	T_1	T_2	T_3	T_4
'	T^{*1+}	T_1^{*1+}	T_2^{*1+}	T_3^{*1+}	T_4^{*1+}
	T^{*+}	T_1^{*+}	T_2^{*+}	T_3^{*+}	
%	$T^{* \%+}$	$T_1^{* \%+}$	$T_2^{* \%+}$		
&	$T^{* \&+}$	$T_1^{* \&+}$			
%'	$T^{* \%'+}$				

$$T(m) = \frac{4 T - (2m) T - (m)}{4 - 1}, \quad k = 1, 2, 3, \dots, \quad (7.21)$$

will approximate $\int_a^b f(x)dx$ to accuracy $O(h^k)$, provided of course that f exists and is continuous on the closed interval $[a, b]$. This extrapolation process is known as the **Romberg** integration method.

The intermediate results in Romberg's method are often arranged in the form of a table, known as the Romberg table. For example, if we start with $m = 4$ subdivisions of the closed interval $[a, b]$, each of length $h = (b - a)/4$, and proceed by doubling the number of subdivisions in each step (and thereby halving the spacing h between the quadrature points from the previous step), then the associated Romberg table is as shown in Table 7.2, where we took, successively, $m = 4, 8, 16, 32, 64$ subdivisions of the interval $[a, b]$ of length $h = (b - a)/m$ each. After $T(4) = T_1(4), \dots, T(64) = T_4(64)$ have been computed, we calculate $T(8), \dots, T(32)$ using (7.21) with $k = 1$, then we compute $T(16), \dots, T(16)$ using (7.21) with $k = 2$, then $T(8), T(8)$ using (7.21) with $k = 3$, and finally $T(4)$ using (7.21) with $k = 4$. Provided that the integrand is sufficiently smooth, the numbers in the $T(m)$ column approximate the integral to within an error $O(h^4)$; the numbers in the $T_1(m)$ column to within $O(h^3)$, those in the $T_2(m)$ column to $O(h^2)$, those in the $T_3(m)$ column to $O(h)$, and those in the $T_4(m)$ column to within $O(h^0)$.

1
 5 6 10
 (1, F 64 8
 O \$ * C*? #>>4 <F = (5

An example is shown in Table 7.3. This gives the results of calculating the integral

$$\int_0^1 \frac{e^{-x}}{1+4x} dx$$

by Romberg’s method; first the trapezium rule is used successively with $m = 4, 8, 16, 32$ and 64 equal subdivisions of the interval $[0, 1]$ of length $h = (b - a)/m$ each. There are then four stages of extrapolation: Stage 1 involves computing $T(m)$ for $m = 4, 8, 16, 32$; Stage 2 computes $T(m)$ for $m = 4, 8, 16$; Stage 3 calculates $T(m)$ for $m = 4, 8$; and Stage 4 then computes $T(m)$ for $m = 4$. Not only does the extrapolation give an accurate result, but the consistency of the numerical values in the last two columns gives a good deal of confidence in quoting the result 0.220458 correct to six decimal digits. Note that none of the individual composite trapezium rule calculations in the $T(m)$ column gives a result correct to more than three decimal digits – not even $T(64)$ which uses 64 equal subdivisions of $[0, 1]$.

Table 7.3. Romberg table for the calculation of $\int_0^1 (e^{-x}/(1+4x))dx$.

	$T(m)$	$R_1(m)$	$R_2(m)$	$R_3(m)$
$m=4$	0.220	0.220	0.220	0.220
$m=8$	0.220	0.220	0.220	0.220
$m=16$	0.220	0.220	0.220	0.220
$m=32$	0.220	0.220	0.220	0.220
$m=64$	0.220	0.220	0.220	0.220
Extrapolated				0.220458

The success of Romberg integration is only justified if the integrand f satisfies the hypotheses of the Euler–Maclaurin Theorem. As an illustration of this, Table 7.4 shows the result of the same calculation, but for the integral

$$\int_0^1 \frac{1}{x} dx.$$

The function $1/x$ is not differentiable at $x = 0$, so the required conditions are not satisfied for any extrapolation. The numerical results bear this out; they are quite close to the correct value, $3/4$, but the behaviour of the extrapolation does not give any confidence in the accuracy of the result. In fact the extrapolation has not given much improvement

on T (64). The calculation of integrals involving this sort of singularity requires special methods which we shall not discuss here.

We have reached the end of this chapter, but do not despair: the story about numerical integration rules will continue. In Chapter 10 we shall discuss a class of quadrature formulae, generally referred to as Gaussian quadrature rules, which are distinct from the Newton–Cotes formulae considered here. Before doing so, however, in Chapters 8 and 9 we make a brief excursion into the realm of approximation theory.

Table 7.4. Romberg table for the calculation of $\int_0^1 x \, dx$.

$R_{0,0}$	$R_{1,0}$	$R_{2,0}$	$R_{3,0}$	$R_{4,0}$
0.500000	0.333333	0.250000	0.208333	0.181818
0.375000	0.250000	0.208333	0.181818	0.166667
0.312500	0.208333	0.181818	0.166667	0.157895
0.270833	0.181818	0.166667	0.157895	0.153846
0.246094	0.166667	0.157895	0.153846	0.152043
0.231250	0.157895	0.153846	0.152043	0.151515
0.223958	0.153846	0.152043	0.151515	0.151322
0.220704	0.152043	0.151515	0.151322	0.151250
0.219231	0.151515	0.151322	0.151250	0.151234
0.218750	0.151322	0.151250	0.151234	0.151230
0.218594	0.151250	0.151234	0.151230	0.151229
0.218542	0.151234	0.151230	0.151229	0.151229
0.218539	0.151230	0.151229	0.151229	0.151229
0.218539	0.151229	0.151229	0.151229	0.151229

7.8 Notes

The material presented in this chapter is classical. For further details on the theory and practice of numerical integration, we refer to the following texts:

- 1. J. Stoer and R. S. Field, *Methods of Numerical Integration*, Second Edition, Computer Science and Applied Mathematics, Academic Press, Orlando, FL, 1984;
- 2. I. M. Abramowitz and I. Stegun, *Approximate Calculation of Integrals*, translated from Russian by Arthur H. Stroud, ACM Monograph Series, Macmillan, New York, 1962;
- 3. G. M. Phillips, *Numerical Quadrature and Cubature*, Computational Mathematics and Applications, Academic Press, London, 1980.

The first of these is a standard text and contains a huge bibliography of more than 1500 entries. Concerning the implementation of numerical integration rules into mathematical software, the reader is referred to

- 4. J. F. Fournier, *Computational Integration*, SIAM, Philadelphia, 1998.

It includes a comprehensive overview of computational integration techniques based on both numerical and symbolical methods, and an exposition of some more recent number-theoretical, pseudorandom and lattice algorithms; these topics are beyond the scope of the present text.

Exercises

- 7.1 With the usual notation for the Newton–Cotes quadrature formula and using the equally spaced quadrature points $x_k = a + kh$ for $k = 0, 1, \dots, n$ and $n \geq 1$, show that $w_k = w_{n-k}$ for $k = 0, 1, \dots, n$.
- 7.2 By considering the polynomial $[x - (a+b)/2]^n$, $n \geq 1$, and the result of Exercise 1, or otherwise, show that the Newton–Cotes formula using $n + 1$ points x_k , $k = 0, 1, \dots, n$, is exact for all polynomials of degree $n + 1$ whenever n is even.
- 7.3 A quadrature formula on the interval $[-1, 1]$ uses the quadrature points $x_0 = -1$ and $x_n = 1$, where $0 < n < \infty$:

$$\int_{-1}^1 f(x) dx \approx w_0 f(-1) + w_n f(1).$$

The formula is required to be exact whenever f is a polynomial of degree 1. Show that $w_0 = w_n = 1$, independent of the value of n . Show also that there is one particular value of n for which the formula is exact also for all polynomials of degree 2. Find this n , and show that, for this value, the formula is also exact for all polynomials of degree 3.

- 7.4 The Newton–Cotes formula with $n = 3$ on the interval $[-1, 1]$ is

$$\int_{-1}^1 f(x) dx \approx w_0 f(-1) + w_1 f(-1/3) + w_2 f(1/3) + w_3 f(1).$$

Using the fact that this formula is to be exact for all polynomials of degree 3, or otherwise, show that

$$\begin{aligned} 2w_1 + 2w_2 &= 2, \\ 2w_1 + -w_2 &= -1, \end{aligned}$$

and hence find the values of the weights w_0 , w_1 , w_2 and w_3 .

- 7.5 For each of the functions $1, x, x^2, \dots, x^4$, find the difference between $\int_{-1}^1 f(x) dx$ and (i) Simpson's rule, (ii) the formula derived in Exercise 4.

Deduce that for every polynomial of degree 5 formula (ii) is

more accurate than formula (i). Find a polynomial of degree 6 for which formula (i) is more accurate than formula (ii).

- 7.6 Write down the errors in the approximation of $\int_a^b x^j dx$ and $\int_a^b x^j dx$

by the trapezium rule and Simpson's rule. Hence find the value of the constant C for which the trapezium rule gives the correct result for the calculation of

$$\int_a^b (x^2 - Cx) dx,$$

and show that the trapezium rule gives a more accurate result than Simpson's rule when $-\frac{1}{2} < C < \frac{1}{2}$.

- 7.7 Determine the values of $c_j, j = 0, 1, 2$, such that the quadrature rule

$$Q(f) = c_0 f(-1) + c_1 f(0) + c_2 f(1) + c_3 f(2)$$

gives the correct value for the integral $\int_{-1}^2 f(x) dx$

$$\int_{-1}^2 f(x) dx$$

when f is any polynomial of degree 3. Show that, with these values of the weights c_j , and under appropriate conditions on the function f ,

$$\left| \int_{-1}^2 f(x) dx - Q(f) \right| \leq M.$$

Give suitable conditions for the validity of this bound, and a definition of the quantity M .

- 7.8 Writing $T(m)$ for the composite trapezium rule defined in (7.15) and $S(2m)$ for the composite Simpson's rule defined in (7.17), show that

$$S(2m) = -T(2m) - T(m).$$

- 7.9 Suppose that the function f has a continuous fourth derivative on the interval $[a, b]$, and that $T(m)$ denotes the composite trapezium rule approximation to $\int_a^b f(x) dx$, using m subintervals. Show that

$$\frac{T(m) - T(2m)}{T(2m) - T(4m)} \rightarrow \frac{1}{4} \text{ as } m \rightarrow \infty.$$

Using the information in Table 7.3 evaluate this expression for $m = 4, 8, 16$.

- 7.10 With the same notation as in Exercise 9, suppose that the fourth derivative of f is not continuous on $[a, b]$, but that

$$\int_a^b f(x) dx - T(m) = A/m^2 + E(m),$$

where $A > 0$ and A are constants and $\lim_{m \rightarrow \infty} m^2 E(m) = 0$. Determine

$$\lim_{m \rightarrow \infty} \frac{T(m)}{T(2m)} - \frac{T(2m)}{T(4m)}.$$

Suggest a value of ϵ which is consistent with the values of $T(m)$ given in Table 7.4.

- 7.11 The function f has a continuous fourth derivative on the interval $[-1, 1]$. Construct the Hermite interpolation polynomial of degree 3 for f using the interpolation points $x = -1$ and $x = 1$. Deduce that

$$\int_{-1}^1 f(x) dx - [f(-1) + f(1)] = -[f'(-1) - f'(1)] + E,$$

where

$$E = \max_{-1 \leq x \leq 1} |f^{(4)}(x)|.$$

- 7.12 Construct the polynomials q_0, q_1, q_2 and q_3 given by Definition 7.3. Hence show that, in the notation of Theorem 7.4,

$$c_0 = 1/12, \quad c_1 = 1/720, \quad c_2 = 1/30240.$$

- 7.13 Using the relations

$$2 \sin x \cos jx = \sin(j+1)x + \sin(j-1)x,$$

$$2 \sin x \sin jx = \cos(j-1)x - \cos(j+1)x,$$

where m is a positive integer, show that the composite trapezium rule (7.15) with m subintervals will give the exact result for each of the integrals

$$\int_0^{\pi} \cos rx \, dx, \quad \int_0^{\pi} \sin rx \, dx,$$

for any integer value of r which is not a multiple of m .

What values are given by the composite trapezium rule for these integrals when $r = mk$ and k is a positive integer?

Polynomial approximation in the ∞ -norm

8.1 Introduction

In Chapter 6 we considered the problem of interpolating a function by polynomials of a certain degree. Here we shall discuss other types of approximation by polynomials, the overall objective being to find the polynomial of given degree n which provides the ‘best approximation’ from \mathcal{P}_n to a given function in a sense that will be made precise below.

8.2 Normed linear spaces

In order to be able to talk about ‘best approximation’ in a rigorous manner we need to recall from Chapter 2 the concept of *norm*; this will allow us to compare various approximations quantitatively and select the one which has the smallest approximation error. The definition given in Section 2.7 applies to a linear space consisting of functions in the same way as to the finite-dimensional linear spaces considered in Chapter 2.

Definition 8.1 Suppose that V is a linear space over the field F of real numbers. A nonnegative function $\| \cdot \|$ defined on V whose value at $f \in V$ is denoted by $\|f\|$ is called a **norm** on V if it satisfies the following axioms:

- $\|f\| = 0$ if, and only if, $f = 0$ in V ;
- $\|cf\| = |c| \|f\|$ for all $c \in F$, and all f in V ;
- $\|f+g\| \leq \|f\| + \|g\|$ for all f and g in V (the triangle inequality).

A linear space V , equipped with a norm, is called a **normed linear space**.

Throughout this chapter [a, b

Lemma 8.1 (i) Suppose that the real-valued weight function w is defined, continuous, positive and integrable on the interval (a, b) . Then, for any function $f \in C[a, b]$,

$$\|f\|_W = \int_a^b w(x) |f(x)| dx, \quad \text{where } W = \int_a^b w(x) dx.$$

(ii) Given any two positive numbers ϵ (however small) and M (however large), there exists a function $f \in C[a, b]$ such that

$$\|f\|_W < \epsilon, \quad \|f\|_\infty > M.$$

Proof The proof is left as an exercise (see Exercise 1). □

The definitions (2.33) and (2.34) of the vector norms $\|\cdot\|_1$ and $\|\cdot\|_2$ on \mathbb{R}^n imply that

$$\|x\|_1 \leq \|x\|_2 \leq \sqrt{2} \|x\|_1, \tag{8.3}$$

which means that, to all intents and purposes, these two norms are interchangeable. Lemma 8.1 indicates that a similar chain of inequalities cannot possibly hold for the norms (8.1) and (8.2) on $C[a, b]$, and the choice between them may therefore significantly influence the outcome of the analysis.

Stimulated by the first axiom of *norm*, we shall think of $f \in C[a, b]$ as being well approximated by a polynomial p on $[a, b]$ if $\|f - p\|$ is small, where $\|\cdot\|$ is either $\|\cdot\|_1$ or $\|\cdot\|_2$ defined, respectively, by (8.1) or (8.2). In the light of Lemma 8.1, it should come as no surprise that the mathematical tools for the analysis of smallness of $\|f - p\|$ are quite different from those that ensure smallness of $\|f - p\|_\infty$. We have therefore chosen to discuss these two matters separately: the present chapter focuses on the ∞ -norm (8.1), while Chapter 9 explores the use of the 2-norm (8.2).

Despite the fundamental differences between the norms (8.1) and (8.2) which we have alluded to above, there is a common underlying feature which is independent of the choice of norm: if *no limitation is imposed*

on the degree of the approximating polynomial p , then the approximation error $\|f - p\|$ can be made *arbitrarily small* in both norms. This is a central result in the theory of polynomial approximation and is formulated in the next theorem.

Theorem 8.1 (Weierstrass Approximation Theorem) *Suppose that f is a real-valued function, defined and continuous on a bounded closed interval $[a, b]$ of the real line; then, given any $\epsilon > 0$, there exists a polynomial p such that*

$$\|f - p\| < \epsilon.$$

Further, if w is a real-valued function, defined, continuous, positive and integrable on (a, b) , then an analogous result holds in the 2-norm over the interval $[a, b]$ with weight function w .

This is an important theorem in classical analysis, and several proofs are known. It is evidently sufficient to consider only the interval $[0, 1]$; a simple change of variable will then extend the proof to any bounded closed interval $[a, b]$. For a real-valued function f , defined and continuous on the interval $[0, 1]$, Bernstein's proof uses the polynomial

$$p_n(x) = \sum_{k=0}^n \binom{n}{k} x^k (1-x)^{n-k} f(k/n), \quad x \in [0, 1],$$

where the **Bernstein polynomials** $p_n(x)$ are defined by

$$p_n(x) = \sum_{k=0}^n \binom{n}{k} x^k (1-x)^{n-k}, \quad x \in [0, 1].$$

It can then be shown that, for any $\epsilon > 0$, there exists $n = n(\epsilon)$ such that $\|f - p_n\| < \epsilon$. The second part of the theorem is a direct consequence of this result, using part (i) of Lemma 8.1.

The details of the proof are given in Exercise 12. For an alternative proof, the reader is referred to Theorem 6.3 in M.J.D. Powell, *Approximation Theory and Methods*, Cambridge University Press, 1996.

1
6 C #@ 9 6 #.@&
1
6 1 1 , 1
1 4 " 15
O + E 5 5 " 6 : ' " "
' \$: V , + +4

B*# H, #.> H 1 5 F \$
F 6D 1 G 6 1 1 \$
6 4 O 1
, 1 , , 5 ES
6 1 G
9 F OI
E 5 5 " 6 : ' " "

8.3 Best approximation in the ∞ -norm

According to the Weierstrass Approximation Theorem any function f in $C[a, b]$ can be approximated arbitrarily well from the set of *all* polynomials. Clearly, if instead of the set of all polynomials we restrict ourselves to the set of polynomials \mathcal{P}_n of degree n or less, with n *fixed*, then it is no longer true that, for any $f \in C[a, b]$ and any $\epsilon > 0$, there exists $p \in \mathcal{P}_n$ such that

$$\|f - p\|_\infty < \epsilon.$$

Consider, for example, the function $x \mapsto \sin x$ defined on the interval $[0, \pi]$ and fix $n = 0$; then $\|f - q\|_\infty \geq 1/2$ for any $q \in \mathcal{P}_0$, and therefore there is no q in \mathcal{P}_0 such that $\|f - q\|_\infty < 1/2$. A similar situation will arise if \mathcal{P}_n is replaced by \mathcal{P}_m , with the polynomial degree n fixed.

It is therefore relevant to enquire just how well a given function f in $C[a, b]$ may be approximated by polynomials of a fixed degree $n \geq 0$. This question leads us to the following approximation problem.

(A) Given that $f \in C[a, b]$ and $n \geq 0$, fixed, find $p \in \mathcal{P}_n$ such that

$$\|f - p\|_\infty = \inf_{q \in \mathcal{P}_n} \|f - q\|_\infty;$$

such a polynomial p is called a **polynomial of best approximation of degree n to the function f in the ∞ -norm**.

The next theorem establishes the existence of a polynomial of best approximation, showing, in particular, that the infimum of $\|f - q\|_\infty$ over $q \in \mathcal{P}_n$ is attained. We shall consider the question of uniqueness of the polynomial of best approximation later on, in Theorem 8.5.

Theorem 8.2 *Given that $f \in C[a, b]$, there exists a polynomial $p \in \mathcal{P}_n$ such that $\|f - p\|_\infty = \min_{q \in \mathcal{P}_n} \|f - q\|_\infty$.*

Proof Let us define the function $E(c_0, \dots, c_n)$ of $n + 1$ real variables by

$$E(c_0, \dots, c_n) = \|f - q\|_\infty, \text{ where } q(x) = c_0 + c_1 x + \dots + c_n x^n.$$

$$\|f - q\|_\infty = \max_{x \in [a, b]} |f(x) - q(x)| = \max_{x \in [a, b]} |f(x) - (c_0 + c_1 x + \dots + c_n x^n)|$$

$$= \max_{x \in [a, b]} |f(x) - c_0 - c_1 x - \dots - c_n x^n| = \max_{x \in [a, b]} |f(x) - c_0 - c_1 x - \dots - c_n x^n|$$

We shall first show that E is continuous; this will imply that E attains its bounds on any bounded closed set in C . We shall then construct a nonempty bounded closed set \mathcal{V} such that the lower bound of E on \mathcal{V} is the same as its lower bound over the whole of C .

To show that E is continuous at each point (c_0, \dots, c_n) , consider any (x_0, \dots, x_n) and define the polynomial $q(x) = \sum_{i=0}^n x_i x^i$ & by $f(q) = \max_{x \in [-1, 1]} |f(x) - q(x)|$. We see from the triangle inequality that

$$\begin{aligned} E(c_0 + \delta_0, \dots, c_n + \delta_n) &= f(q + \delta) \\ &= f(q) + \delta \\ &= E(c_0, \dots, c_n) + \delta. \end{aligned}$$

Now, for any given positive number ϵ , choose $\delta = \epsilon / (1 + K)$, where $K = \max_{i=0, \dots, n} |a_i|$. Consider any (x_0, \dots, x_n) such that $|x_i - c_i| < \delta$ for all $i = 0, \dots, n$. Then,

$$\begin{aligned} E(c_0 + \delta_0, \dots, c_n + \delta_n) &\leq E(c_0, \dots, c_n) \\ &\quad + \max_{i=0, \dots, n} (|a_i| \delta + \delta) \\ &= (1 + K) \delta \\ &= \epsilon. \end{aligned} \tag{8.4}$$

Similarly,

$$\begin{aligned} E(c_0, \dots, c_n) &= f(q) = f(q + \delta) + \\ &\quad - f(q + \delta) + \\ &= E(c_0 + \delta_0, \dots, c_n + \delta_n) - \delta, \end{aligned}$$

and therefore

$$E(c_0, \dots, c_n) \leq E(c_0 + \delta_0, \dots, c_n + \delta_n). \tag{8.5}$$

From (8.4) and (8.5) we deduce that

$$E(c_0 + \delta_0, \dots, c_n + \delta_n) \leq E(c_0, \dots, c_n)$$

for all (x_0, \dots, x_n) such that $|x_i - c_i| < \delta$, $i = 0, \dots, n$, where now $\delta = \epsilon / (1 + K)$ and $K = \max_{i=0, \dots, n} |a_i|$. Hence E is continuous at (c_0, \dots, c_n) . Since (c_0, \dots, c_n) is an arbitrary point in C , it follows that E is continuous on the whole of C .

Let us denote by \mathcal{V} the set of all points (c_0, \dots, c_n) in C such that $E(c_0, \dots, c_n) \leq f + 1$. The set \mathcal{V} is evidently bounded and closed in C ; further, \mathcal{V} is nonempty since $E(0, \dots, 0) = f + 1$, so that $(0, \dots, 0) \in \mathcal{V}$. Hence the continuous function E attains its

lower bound over the set S ; let us denote this lower bound by d and let (c_0, \dots, c_n) denote the point in S where it is attained.

Since $(0, \dots, 0) \in S$, it follows that

$$d = \min_{(c_0, \dots, c_n) \in S} E(c_0, \dots, c_n) - E(0, \dots, 0) = f - f_0.$$

According to the definition of d ,

$$E(c_0, \dots, c_n) - E(0, \dots, 0) \geq d.$$

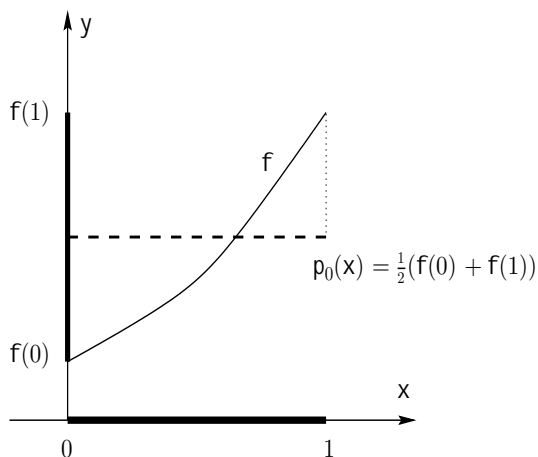


Figure 8.1

The polynomial p will be of the form $p(x) = c$, and we need to determine c so that

$$E(c) = \max_{x \in [0, 1]} |f(x) - c|$$

is minimal. Since f is monotonic increasing, $f(x) - c$ attains its minimum at $x = 0$ and its maximum at $x = 1$; therefore $|f(x) - c|$ reaches its maximum value at one of the endpoints of $[0, 1]$, i.e.,

$$E(c) = \max_{x \in [0, 1]} |f(x) - c| = \max\{|f(0) - c|, |f(1) - c|\}.$$

Clearly,

$$E(c) = \begin{cases} f(1) - c & \text{if } c < \frac{f(0) + f(1)}{2}, \\ c - f(0) & \text{if } c \geq \frac{f(0) + f(1)}{2}. \end{cases}$$

Drawing the graph of the function $c \mapsto E(c)$ shows that the minimum is attained when $c = \frac{f(0) + f(1)}{2}$. Consequently, the desired minimax polynomial of degree 0 for the function f is

$$p(x) = \frac{f(0) + f(1)}{2}, \quad x \in [0, 1].$$

The function f and its minimax approximation p are depicted in Figure 8.1.

More generally, if $f \in C[a, b]$ (not necessarily monotonic), and α and β denote two points in $[a, b]$ where f attains its minimum and maximum

values, respectively, then the minimax polynomial of degree 0 to f on $[a, b]$ is

$$p(x) = \frac{1}{2}(f(a) + f(b)), \quad x \in [a, b].$$

This example shows that the minimax polynomial p of degree zero for $f \in C[a, b]$ has the property that the approximation error $f - p$ attains its extrema at *two* points, $x = a$ and $x = b$, with the error

$$f(x) - p(x) = \frac{1}{2}(f(x) - f(a)) + \frac{1}{2}(f(x) - f(b))$$

being *negative at one point*, $x = a$, and *positive at the other*, $x = b$. We shall prove that a property of this kind holds in general; the precise formulation of the general result is given in Theorem 8.4 which is, due to the oscillating nature of the approximation error, usually referred to as the Oscillation Theorem: it gives a complete characterisation of the minimax polynomial and provides a method for its construction. We begin with a preliminary result due to de la Vallée Poussin.

Theorem 8.3 (De la Vallée Poussin’s Theorem) *Let $f \in C[a, b]$ and $r \in C[a, b]$. Suppose that there exist $n + 2$ points $x_0 < x_1 < \dots < x_{n+1}$ in the interval $[a, b]$, such that $f(x_i) - r(x_i) = (-1)^i \mu$ and $f(x_{i+1}) - r(x_{i+1}) = (-1)^{i+1} \mu$ have opposite signs, for $i = 0, \dots, n$. Then,*

$$\min_p \max_{x \in [a, b]} |f(x) - p(x)| = \mu. \tag{8.6}$$

Proof The condition on the signs of $f(x_i) - r(x_i)$ is usually expressed by saying that $f - r$ has alternating signs at the points x_i , $i = 0, 1, \dots, n+1$. Let us denote the right-hand side of (8.6) by μ . Clearly, $\mu \geq 0$; when $\mu = 0$ the statement of the theorem is trivially true, so we shall assume that $\mu > 0$. Suppose that (8.6) is false; then, for a minimax polynomial approximation p to the function f we have

$$\max_{x \in [a, b]} |f(x) - p(x)| < \mu.$$

¹ K. F. 5, (- B# % #.?? : \$
 5 C , #@? : 5 D #.@ , 1
 8 6 6 5 4
² , 1 .4 , 8 4

Therefore,

$$p(x) - f(x) < r(x) - f(x), \quad i = 0, 1, \dots, n+1.$$

Now,

$$r(x) - p(x) = [r(x) - f(x)] - [p(x) - f(x)], \quad i = 0, 1, \dots, n+1.$$

Since the first term on the right always exceeds the second term in absolute value, it follows that $r(x) - p(x)$ and $r(x) - f(x)$ have the same sign for $i = 0, 1, \dots, n+1$. Hence $r - p$, which is a polynomial of degree n , changes sign $n+1$ times. Thus, the assumption that (8.6) is false has led to a contradiction, and the proof is complete. \square

Theorem 8.3 gives a clue to formulating a constructive characterisation of the *minimax polynomial*: indeed, we shall show that if the quantities $f(x) - r(x)$, $i = 0, 1, \dots, n+1$, in Theorem 8.3 are all equal to $\pm r - f$, then $r - f$ is, in fact, a minimax polynomial of degree n for the function f on the interval $[a, b]$.

Theorem 8.4 (The Oscillation Theorem) *Suppose that $f \in C[a, b]$. A polynomial $r - f$ is a minimax polynomial for f on $[a, b]$ if, and only if, there exists a sequence of $n+2$ points x_i , $i = 0, 1, \dots, n+1$, such that $a < x_0 < x_1 < \dots < x_{n+1} < b$,*

$$f(x_i) - r(x_i) = (-1)^i (r - f), \quad i = 0, 1, \dots, n+1,$$

and

$$f(x) - r(x) = (-1)^i (f(x_i) - r(x_i)), \quad i = 0, \dots, n.$$

The statement of the theorem is often expressed by saying that $f - r$ attains its maximum absolute value with alternating signs at the points x_i . The points x_i , $i = 0, 1, \dots, n+1$, in the Oscillation Theorem are referred to as **critical points**.

Proof of theorem If $f - r$ is a minimax polynomial, then the result is trivially true, with $r = f$ and any sequence of $n+2$ distinct points x_i , $i = 0, 1, \dots, n+1$, contained in $[a, b]$. Thus, we shall suppose throughout the proof that $f - r \neq 0$, i.e., f is such that there is no polynomial $p - f$ whose restriction to $[a, b]$ is identically equal to f .

The sufficiency of the condition stated in the theorem is easily shown. Suppose that the sequence of points x_i , $i = 0, 1, \dots, n+1$, exists with

the given properties. Define

$$L = \|f - r\| \quad \text{and} \quad E(f) = \min_P \|f - q\|.$$

From De la Vallée Poussin's Theorem, Theorem 8.3, it follows that $E(f) \leq L$. By the definition of $E(f)$ we also see that $E(f) \leq \|f - r\| = L$. Hence $E(f) = L$, and the given polynomial r is a minimax polynomial.

For the necessity of the condition, suppose that the given polynomial r & is a minimax polynomial for f on $[a, b]$. As $x \mapsto f(x) - r(x)$ is a continuous function on the bounded closed interval $[a, b]$, there exists a point in $[a, b]$ at which $f(x) - r(x)$ attains its maximum value, $L > 0$; let

$$x = \min_{x \in [a, b]} \{f(x) - r(x)\} = L.$$

Now, $x = b$ would imply that $f(x) - r(x) = L$ for all $x \in [a, b]$. As f is continuous on $[a, b]$, it would then follow that either $f(x) = r(x) + L$ for all $x \in [a, b]$ or $f(x) = r(x) - L$ for all $x \in [a, b]$; either way, we would find that f & , which is assumed not to be the case. Therefore, $x \in [a, b]$; we may assume without loss of generality that $f(x) - r(x) = L > 0$.

Now, we shall prove the existence of the next critical point, $x \in (x, b]$ such that $f(x) - r(x) = -L$. Suppose otherwise, for contradiction; then, $-L < f(x) - r(x) < L$ for all x in $[a, b]$. Thus, by the continuity of f , there exists $(0, L)$ such that $L + \epsilon = f(x) - r(x) < L$ for all $x \in [a, b]$. Let us define r & by

$$r(x) = r(x) + \epsilon,$$

where $0 < \epsilon < \min\{\epsilon, L - \epsilon\}$. Then, for all $x \in [a, b]$,

$$f(x) - r(x) = f(x) - r(x) - \epsilon < L + \epsilon < L$$

and

$$f(x) - r(x) = f(x) - r(x) - \epsilon < L - \epsilon < L,$$

which means that

$$\|f - r\| < L = \|f - r\|.$$

Hence, r & is a better approximation to f on $[a, b]$ than r & is. This, however, contradicts our hypothesis that r is a polynomial of best approximation to f on $[a, b]$ from \mathcal{P}_n , and implies the existence of

$$x = \inf_{x \in (x, b]} \{f(x) - r(x)\} = -L.$$

Consequently, $f(x) - r(x) = 0$ and $x \in (x, b]$, as required; thus if $n = 0$, the proof is complete.

Let us, therefore, suppose that n

on each of the intervals $[x_i, x_{i+1}]$, $i = 0, 1, \dots, m$ (whose union is $[a, b]$). We shall prove that, for $\delta > 0$ sufficiently small,

$$f(x) - r(x) < L = f - r$$

for all x in $[x_i, x_{i+1}]$ and all $i = 0, 1, \dots, m$; i.e., $f - r < f - r$, contradicting the fact that r & is a minimax polynomial for f on $[a, b]$, and refuting the hypothesis that $1 - m - n$.

Take, for example, the interval $[x_i, x_{i+1}]$. For each x in $[x_i, x_{i+1}]$ we have $v(x) > 0$ and therefore, by the definition of $r(x)$ and property (a) above,

$$f(x) - r(x) - L - v(x) < L, \quad x \in [x_i, x_{i+1}].$$

Further, as $v(x_i) = 0$, it follows from (d) that

$$f(x_i) - r(x_i) = f(x_i) - r(x_i) < L.$$

Therefore, $f(x) - r(x) < L$ for each x in $[x_i, x_{i+1}]$. For a lower bound on $f(x) - r(x)$, note that by (a) and (c), $f(x) - r(x) > L$ for all x in $[x_i, x_{i+1}]$. As $f - r$ is a continuous function on $[x_i, x_{i+1}]$, there exists $(0, L)$ such that $f(x) - r(x) > L + \delta$ for all x in $[x_i, x_{i+1}]$. Thus, for $0 < \delta < \min L, \delta, \delta$, where

$$\delta = \frac{1}{\max_{x \in [x_i, x_{i+1}]} v(x)},$$

we have that

$$f(x) - r(x) > L + \delta - v(x) > L, \quad x \in [x_i, x_{i+1}].$$

Further, by (d) above,

$$f(x_i) - r(x_i) = f(x_i) - r(x_i) > L.$$

Hence, $f(x) - r(x) > L$ for all $x \in [x_i, x_{i+1}]$, for $0 < \delta < \min L, \delta, \delta$. Combining the upper and lower bounds on $f(x) - r(x)$, we deduce that

$$f(x) - r(x) < L = f - r, \quad x \in [x_i, x_{i+1}].$$

Arguing in the same manner on each of the other intervals $[x_i, x_{i+1}]$, $i = 1, \dots, m$, with $0 < \delta < \min L, \delta, \delta$, $i = 1, \dots, m$, and δ defined analogously to δ and δ above, we conclude that

$$f(x) - r(x) < L = f - r, \quad x \in [x_i, x_{i+1}], \quad i = 0, 1, \dots, m,$$

and hence, for $0 < \delta < \min L, \delta, \delta, \delta, \delta, \delta$,

$$f - r < L = f - r.$$

of replacing r by $r(x) = r(x) + v(x)$, with $v > 0$, is indicated by the arrows. Since $f - r = f - r + v(x)$ and v is negative for $x \in (c, d)$ and positive outside (c, d) , $f - r$ will be smaller than $f - r$ at each of the points P_i , $i = 0, 1, 2$. There are two other local extrema for the error function $f - r$: a minimum at Q and a maximum at R . Since both these points are to the right of c , where $v(x) > 0$, we shall have $f - r > f - r$ at both of Q and R , and $f - r > f - r$ at R . The magnitude of the extra term $v(x)$ must therefore be limited by the need to avoid the new difference $f - r$ becoming too large at R . We can achieve this by selecting $v > 0$ sufficiently small. In this illustration the polynomial r & is not a minimax approximation to f on the given interval, since we can construct a better approximation r which is also in & .

We can now apply the Oscillation Theorem to prove that the minimax polynomial is unique.

Theorem 8.5 (Uniqueness Theorem) *Suppose that $[a, b]$ is a bounded closed interval of the real line. Each $f \in C[a, b]$ has a unique minimax polynomial p & on $[a, b]$.*

Proof Suppose that q & is also a minimax polynomial for f , and that p and q are distinct. Then,

$$f - p = f - q = E(f),$$

where, as in the proof of the Oscillation Theorem, we have used the notation

$$E(f) = \min_p f - p.$$

This implies, by the triangle inequality, that

$$\begin{aligned} f - (p + q) &= -(f - p) + -(f - q) \\ &= -E(f) + -E(f) \\ &= -2E(f). \end{aligned}$$

Therefore $-(p + q)$ & is also a minimax polynomial approximation to f on $[a, b]$. By the Oscillation Theorem there exists a sequence of $n + 2$ critical points x_i , $i = 0, 1, \dots, n + 1$, at which

$$f(x_i) - (p(x_i) + q(x_i)) = -2E(f), \quad i = 0, 1, \dots, n + 1.$$

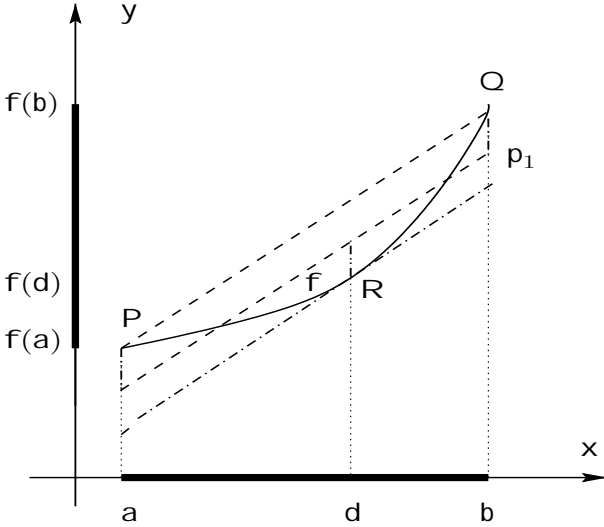


FIGURE 8.3

equations

$$\begin{aligned} f(a) - (c/a + c) &= A, \\ f(d) - (c/d + c) &= A, \\ f(b) - (c/b + c) &= A, \end{aligned} \tag{8.7}$$

where either $A = L$ or $A = -L$, with $L = \max_{a \leq x \leq b} |f(x) - p(x)|$. Along with the condition

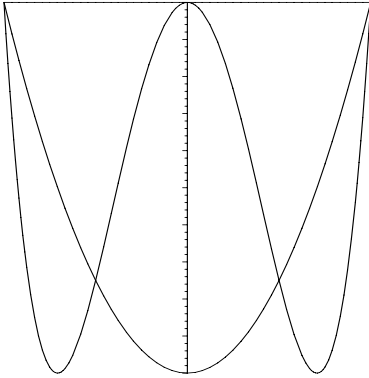
$$f'(d) = c \tag{8.8}$$

this gives four equations to determine the unknowns d, c, c and A .

Subtracting the first equation in (8.7) from the third equation, we get $f(b) - f(a) = c(b - a)$, whereby $c = (f(b) - f(a))/(b - a)$. Now, by the Mean Value Theorem, Theorem A.3, with this choice of c equation (8.8) has at least one solution, d , in the open interval in (a, b) . In fact, the value of d is uniquely determined by (8.8), as f' is continuous and strictly monotonic increasing. Next, c can be determined by adding the second equation in (8.7) to the first. Having calculated both c and c we insert them into the first equation in (8.7) to obtain A ; finally $L = A$.

The construction of the minimax polynomial p_1 is illustrated in Figure 8.3; R is the point at which the tangent to the curve $y = f(x)$ is parallel to the chord PQ ; the graph of $p_1(x)$ is parallel to these two lines, and lies half-way between them.

Table 8.1.



they are all real and distinct, and lie in $(-1, 1)$;

- (v) $T_n(x) \leq 1$ for all $x \in [-1, 1]$ and all $n \geq 0$;
 (vi) for $n \geq 1$, $T_n(x) = \pm 1$, alternately at the $n + 1$ points $x_k = \cos(k\pi/n)$, $k = 0, 1, \dots, n$.

We can now apply the Oscillation Theorem to construct the minimax polynomial of degree n for $f: x \mapsto x$ on the interval $[-1, 1]$.

Theorem 8.6 Suppose that $n \geq 0$. The polynomial p_n is defined by

$$p_n(x) = x - 2^{-n} T_{n+1}(x), \quad x \in [-1, 1],$$

is the minimax approximation of degree n to the function $x \mapsto x$ on the interval $[-1, 1]$.

Proof By part (ii) of Lemma 8.2, p_n is defined by

$$x - p_n(x) = 2^{-n} T_{n+1}(x),$$

by parts (v) and (vi) of Lemma 8.2, the difference $x - p_n(x)$ does not exceed 2^{-n} in the interval $[-1, 1]$, and attains this value with alternating signs at the $n + 2$ points $x_k = \cos(k\pi/(n + 1))$, $k = 0, 1, \dots, n + 1$. Therefore, by the Oscillation Theorem, p_n is the (unique) minimax polynomial approximation from \mathcal{P}_n to the function $x \mapsto x$ over $[-1, 1]$. \square

A polynomial of degree n whose leading coefficient, the coefficient of x^n , is equal to 1, is called a **monic polynomial** of degree n . For example, the polynomial r_{n+1} is defined by $r_{n+1}(x) = x^{n+1} - q_n(x)$ with $q_n \in \mathcal{P}_n$, is a monic polynomial of degree $n + 1$.

Corollary 8.1 Suppose that $n \geq 0$. Among all monic polynomials of degree $n + 1$ the polynomials $2^{-n} T_{n+1}$ and $-2^{-n} T_{n+1}$ have the smallest ∞ -norm on the interval $[-1, 1]$.

Proof Let \mathcal{P}_{n+1} denote the set of all monic polynomials of degree $n + 1$. Any $r \in \mathcal{P}_{n+1}$ can be regarded as the difference between the function $x \mapsto x^{n+1}$ and a polynomial of lower degree, i.e., $r(x) = x^{n+1} - q(x)$ with $q \in \mathcal{P}_n$. Hence, by Theorem 8.6,

$$\begin{aligned} \min_{\mathcal{P}_{n+1}} \|r\|_{\infty} &= \min_{\mathcal{P}_n} \|x^{n+1} - q\|_{\infty} \\ &= \min_{\mathcal{P}_n} \|x^{n+1} - 2^{-n} T_{n+1}\|_{\infty} \\ &= 2^{-n} \|T_{n+1}\|_{\infty}; \end{aligned}$$

the minimum is, therefore, achieved when $r = 0$ &

Proof Let $t_j = \cos(j\pi/(n+1))$, $j = 0, 1, \dots, n$, denote the zeros of the polynomial $T_{n+1}(t)$ (in the interval $(-1, 1)$). Hence,

$$T_{n+1}(t) = 2^{-n} \prod_{j=0}^n (t - t_j).$$

Let us define the points x_j , $j = 0, 1, \dots, n$, as in the statement of the theorem. Clearly (x_j, y_j) is the image of $(t_j, 1)$ under the linear transformation $t = x = -(b-a)t + -(b+a)$; we note in passing that the inverse of this mapping is $x = t(x) = (2x - a - b)/(b - a)$; thus,

$$T_{n+1}(t(x)) = \frac{b-a}{2} \prod_{j=0}^n (t(x) - t_j) = \frac{b-a}{2} 2^{-n} \prod_{j=0}^n (t(x)).$$

The required bound now follows from (8.9), since $|T_{n+1}(t(x))| \leq 1$ for all $x \in [a, b]$, and therefore $|f(x) - p(x)| \leq (b-a) 2^{-n} 2^{-n}$. \square

The De la Vallée Poussin Theorem, Theorem 8.3, suggests the notion of a **near-minimax** polynomial, which is a polynomial p such that the difference $f(x) - p(x)$ changes sign at $n + 1$ points x_j , $j = 0, 1, \dots, n$, with $a < x_0 < x_1 < \dots < x_n < b$; for the difference $f(x) - p(x)$ then attains a local maximum or minimum with alternating signs in each of the intervals $[a, x_0], (x_0, x_1), \dots, (x_n, b]$. The positions of these alternating local maxima and minima are then the points x_i , $i = 0, 1, \dots, n + 1$, required by Theorem 8.3, and we therefore know that the ∞ -norm of the error of the minimax polynomial lies between the least and greatest of the absolute values of these local maxima and minima. In particular, we should expect that if the sizes of these local maxima and minima are not greatly different, then the error of the near-minimax approximation should not be very much larger than the error of the minimax approximation.

Given any set of points x_i , $i = 0, 1, \dots, n$, with $a < x_0 < x_1 < \dots < x_n < b$, the polynomial $\omega(x) = (x - x_0) \dots (x - x_n)$ changes sign at the $n + 1$ points x_j , $j = 0, 1, \dots, n$. Let us assume that $f \in C[a, b]$, f' exists and is continuous on $[a, b]$, and f' has the same sign on the whole of (a, b) . It then follows that the product $f'(x) \omega(x)$ has exactly $n + 1$ sign-changes in the open interval (a, b) for any (a, b) . Thus, according to (8.9), the Lagrange interpolation polynomial p of degree n for the function f , with interpolation points x_j , $j = 0, 1, \dots, n$, contained in the open interval (a, b) , is a near-minimax polynomial from \mathcal{P}_n for f on $[a, b]$.

We have therefore just shown that if f' exists and is continuous on the closed interval $[a, b]$, and has the same sign on the open interval

so we obtain a polynomial approximation $p_n(x)$ by taking the terms of this series up to the one involving x^n . Then, clearly,

$$e^x - p_n(x) = \frac{x^{n+1}}{(n+1)!}.$$

Over the interval $[0, 1]$, for example, this difference is nonnegative and monotonic increasing; it does not change sign at all. Hence the polynomial p_n thus constructed is quite certainly not a near-minimax approximation for $x \in [0, 1]$. Nevertheless, $\max_{x \in [0, 1]} (e^x - p_n(x))$ can be made arbitrarily small by choosing n sufficiently large.

8.6 Notes

For further details on the topics presented in this chapter, we refer to

+ [1] *Approximation Theory and Methods*, Cambridge University Press, Cambridge, 1996.

The Weierstrass Theorem is discussed in Chapter 6 of that book, and is stated in its Theorem 6.3. Although the proof presented by Powell uses the Bernstein polynomials, it is different from the more elementary but slightly lengthier argument proposed in Exercise 12 here: it relies on a proof of Bohman and Korovkin based on properties of monotone operators; see, also, p. 66 in Chapter 3 of

[2] *Introduction to Approximation Theory*, McGraw-Hill, New York, 1966.

The notes contained on pp. 224–233 of Cheney's book are particularly illuminating.

The proof of the Weierstrass Theorem as proposed in Exercise 12, including the definition of what we today call Bernstein polynomials, stem from a paper of Sergei Natanovich Bernstein (1880–1968), entitled 'Démonstration du théorème de Weierstrass fondée sur le calcul des probabilités', *Comm. Soc. Math. Kharkov* **13**, 1–2, 1912/13.

Weierstrass' main contributions to approximation theory, as well as those of other mathematicians (including Picard, Volterra, Runge, Lebesgue, Mittag-Leffler, Fejér, Landau, de la Vallée Poussin, Bernstein), are reviewed in the extensive historical survey by Allan Pinkus, *Weierstrass and approximation theory*, *J. Approx. Theory* **107**, 1–66, 2000. Further details about the history of the subject can be found at

the history of approximation theory website maintained by Allan Pinkus and Carl de Boor: <http://www.math.umd.edu/~pinkus/10>

The second part of Theorem 8.1 concerning the approximability of a continuous function by polynomials in the 2-norm is not usually presented as part of the classical Weierstrass Theorem which is posed in the ∞ -norm. Here, we have chosen to state these results together in order to highlight the analogy, as well as to motivate the use of the 2-norm in polynomial approximation in the next chapter, Chapter 9.

In both Cheney's and Powell's books minimax approximation is treated in the more general framework of Haar systems. An $(n + 1)$ -dimensional linear subspace \mathcal{H} of $C[a, b]$ is said to satisfy the *Haar condition* if, for every nonzero p in \mathcal{H} , the number of roots of the equation $p(x) = 0$ in the interval $[a, b]$ is less than $n + 1$. The concept of Haar system is due to Alfred Haar (1885–1933), *Die Minkowskische Geometrie und die Annäherung an stetige Funktionen*, *Math. Ann.* **78**, 294–311, 1918; this paper contains Haar's Theorem which characterises finite-dimensional Haar systems in spaces of continuous functions. The *Characterisation Theorem*, formulated as Theorem 7.2 in Powell's book, shows that the Oscillation Theorem, Theorem 8.4 of the present chapter, remains valid in a more general setting when the set of polynomials $1, x, \dots, x^n$ is replaced by an $(n + 1)$ -dimensional Haar system of functions contained in $C[a, b]$.

Exercises

- 8.1 Give a proof of Lemma 8.1.
- 8.2 Suppose that the real-valued function f is continuous and even on the interval $[-a, a]$, that is, $f(x) = f(-x)$ for all $x \in [-a, a]$. By using the Uniqueness Theorem, or otherwise, show that the minimax polynomial approximation of degree n is an even function. Deduce that the minimax polynomial approximation of degree $2n$ is also the minimax polynomial approximation of degree $2n + 1$. What does this imply about the sequence of critical points for the minimax polynomial p_n ?
- 8.3 State and prove similar results to those in Exercise 2, for the case where f is an odd function, that is, $f(x) = -f(-x)$ for all $x \in [-a, a]$.
- 8.4 (i) Construct the minimax polynomial p_n & q_n on the interval $[-1, 1]$ for the function g defined by $g(x) = \sin x$.

(ii) Construct the minimax polynomial p_n on the interval $[-1, 1]$ for the function h defined by $h(x) = \cos x$.
 (Use the results of Exercises 2 and 3.)

8.5 The function H is defined by $H(x) = 1$ if $x > 0$, $H(x) = -1$ if $x < 0$, and $H(0) = 0$. Show that for any $n \geq 0$ and any p_n , $\|H - p_n\| \geq 1/2$ on the interval $[-1, 1]$. Construct the polynomial, of degree n , of best approximation to H on the interval $[-1, 1]$, and show that it is unique. (Note that since H is discontinuous most of the theorems in this chapter are not applicable.)

Show that the polynomial of best approximation, of degree n , to H on $[-1, 1]$ is not unique, and give an expression for its most general form.

8.6 Suppose that $t_0 < t_1 < \dots < t_k$ are k distinct points in the interval $[a, b]$; for any function f defined on $[a, b]$, write $Z(f) = \max_{0 \leq i \leq k} |f(t_i)|$. Explain why $Z(\cdot)$ is not a norm on the space of functions which are continuous on $[a, b]$; show that it is a norm on the space of polynomials of degree n , provided that $k > n$.

In the case $k = 3$, with $t_0 = 0$, $t_1 = -1/2$, $t_2 = 1/2$, $t_3 = 1$, where we wish to approximate the function $f: x \mapsto e^x$ on the interval $[0, 1]$, explain graphically, or otherwise, why the polynomial p of degree 1 which minimises $Z(f - p)$ satisfies the conditions

$$f(0) - p(0) = [f(-1/2) - p](-1/2) = f(1/2) - p(1/2).$$

Hence construct this polynomial p . Now suppose that $k = 4$, with $t_0 = 0$, $t_1 = -1/2$, $t_2 = 1/2$, $t_3 = 1$; use a similar method to construct the polynomial of degree 1 which minimises $Z(f - p)$.

8.7 Among all polynomials p_n of the form

$$p_n(x) = Ax^n + a_{n-1}x^{n-1} + \dots + a_0,$$

where A is a fixed nonzero real number, find the polynomial of best approximation for the function $f(x) = \ln x$ on the closed interval $[1/2, 1]$.

8.8 Find the minimax polynomial p_n on the interval $[-1, 1]$ for the function f defined by

$$f(x) = \frac{1}{1+x^2},$$

where $a = 0$.

- 8.9 Construct the minimax polynomial p_n on the interval $[1, 2]$ for the function f defined by $f(x) = x$.
- 8.10 Give a proof of Lemma 8.2.
- 8.11 Give an example of a continuous real-valued function f defined on the closed interval $[a, b]$ such that the set of critical points for the minimax approximation of f by polynomials from \mathcal{P}_n does not contain either of the points a and b .
- 8.12 For each nonnegative integer n , and $x \in [0, 1]$, define the Bernstein polynomials p_n by

$$p_n(x) = \frac{n!}{k!(n-k)!} x^k (1-x)^{n-k}, \quad k = 0, \dots, n.$$

Show that

$$(1-x+tx)^n = \sum_{k=0}^n p_k(x)t^k;$$

by differentiating this relation successively with respect to t and putting $t = 1$, show that, for any $x \in [0, 1]$,

$$p_n(x) = 1,$$

$$k p_k(x) = n x,$$

$$k(k-1) p_k(x) = n(n-1) x^2,$$

and deduce that

$$\left(x - \frac{k}{n}\right) p_k(x) = \frac{x(1-x)}{n}, \quad x \in [0, 1].$$

Define M to be the upper bound of $f(x)$ on $[0, 1]$. Given $\epsilon > 0$, we can choose $\delta > 0$ such that $f(x) - f(y) < \epsilon/2$ for any x and y in $[0, 1]$ such that $|x - y| < \delta$. Now define the polynomial p_n by

$$p_n(x) = \sum_{k=0}^n f(k/n) p_k(x),$$

and choose a fixed value of x in $[0, 1]$; show that

$$f(x) - p_n(x) = \sum_{k=0}^{n-1} (f(x) - f(k/n)) p_k(x).$$

Using the notation

$$S_1 = \sum_{k=0}^{n-1} \dots + \dots$$

where S_1 denotes the sum over those values of k for which $x - k/n < \dots$, and S_2 denotes the sum over those values of k for which $x - k/n \geq \dots$, show that

$$|f(x) - p_n(x)| < \dots / 2.$$

Show also that

$$|f(x) - p_n(x)| \leq (2M/n) \sum_{k=0}^{n-1} p_k(x).$$

Now, choose $N = M/\dots$, and show that

$$|f(x) - p_n(x)| < \dots \quad x \in [0, 1],$$

if $n \geq N$. Deduce that

$$\|f - p_n\| < \dots, \quad \text{if } n \geq N,$$

where $\|\cdot\|$ denotes the ∞ -norm on the interval $[0, 1]$.

Approximation in the 2-norm

9.1 Introduction

In Chapter 8 we discussed the idea of best approximation of a continuous real-valued function by polynomials of some fixed degree in the ∞ -norm. Here we consider the analogous problem of best approximation in the 2-norm. Why, you might ask, is it necessary to consider best approximation in the 2-norm when we have already developed a perfectly adequate theory of best approximation in the ∞ -norm? As our first example in Section 9.3 will demonstrate, the choice of norm can significantly influence the outcome of the problem of best approximation: the polynomial of best approximation of a certain fixed degree to a given continuous function in one norm need not bear any resemblance to the polynomial of best approximation of the same degree in another norm. Ultimately, in a practical situation, the choice of norm will be governed by the sense in which the given continuous function has to be well approximated.

As will become apparent, best approximation in the 2-norm is closely related to the notion of orthogonality and this in turn relies on the concept of *inner product*. Thus, we begin the chapter by recalling from linear algebra the definition of *inner product space*.

Throughout the chapter $[a, b]$ will denote a nonempty, bounded, closed interval of the real line, and (a, b) will signify a nonempty bounded open interval of the real line.

Now,

$$\begin{aligned} \langle u, v \rangle &= \|u\| \|v\| \cos(\theta) \\ &= \|u\| \|v\| (\cos \theta \cos \phi + \sin \theta \sin \phi) \\ &= \|u\| \|v\| \cos(\theta - \phi) \\ &= \|u\| \|v\| \cos(\theta - \phi), \end{aligned}$$

where $\theta = \angle(u, v)$ is the angle between the vectors u and v . The vector u is orthogonal to v if, and only if, $\theta = \pi/2$ or $3\pi/2$; either way, $\cos(\theta - \phi) = 0$, and hence $\langle u, v \rangle = 0$. We note in passing that if $\theta = \phi$, then $\theta - \phi = 0$ and therefore

$$\langle u, v \rangle = \|u\| \|v\| \cos(0) = \|u\| \|v\|.$$

This last observation motivates our next definition.

Definition 9.3 Suppose that V is an inner product space over the field of real numbers, with inner product $\langle \cdot, \cdot \rangle$. For f in V , we define the

$$\|f\| = \sqrt{\langle f, f \rangle}. \tag{9.1}$$

Although our terminology and our notation appear to imply that (9.1) defines a norm on V , this is by no means obvious. In order to show that $\|f\| = \sqrt{\langle f, f \rangle}$ is indeed a norm, we begin with the following result which is a direct generalisation of the Cauchy–Schwarz inequality (2.35) from Chapter 2.

Lemma 9.1 (Cauchy–Schwarz inequality)

$$|\langle f, g \rangle| \leq \|f\| \|g\| \tag{9.2}$$

Proof The proof is analogous to that of (2.35). Recalling the definition of $\|f\|$ from (9.1) and noting the first three axioms of inner product, we find that, for $f, g \in V$,

$$0 \leq \|f + g\|^2 = \|f\|^2 + 2\langle f, g \rangle + \|g\|^2. \tag{9.3}$$

Denoting, for $f, g \in V$ fixed, the quadratic polynomial in α on the right-hand side by $A(\alpha)$, the condition for $A(\alpha)$ to be nonnegative for all α in \mathbb{R} is that $[2\langle f, g \rangle]^2 - 4\|f\|^2\|g\|^2 \leq 0$; this gives the inequality (9.2). \square

Now, putting $\alpha = 1$ in (9.3) and using (9.2) on the right yields

$$\|f + g\|^2 \leq \|f\|^2 + \|g\|^2 + 2\langle f, g \rangle.$$

Consequently, $\| \cdot \|$ obeys the triangle inequality, the third axiom of norm. The first two axioms of norm, namely that

$$\begin{aligned} & \|f\| \geq 0 \text{ for all } f, \text{ and } \|f\| = 0 \text{ if, and only if, } f = 0 \text{ in } V, \text{ and} \\ & \|af\| = |a| \|f\| \text{ for all } a \text{ and all } f, \end{aligned}$$

follow directly from (9.1) and from the last three axioms of inner product stated in Definition 9.1.

We have thus shown the following result.

Theorem 9.1 *An inner product space V over the field F of real numbers, equipped with the induced norm $\| \cdot \|$, is a normed linear space over F .*

We conclude this section with a relevant example of an inner product space, whose induced norm is the 2-norm considered at the beginning of Chapter 8.

Example 9.3 *The set $C[a, b]$ of continuous real-valued functions defined on the closed interval $[a, b]$ is an inner product space with*

$$\langle f, g \rangle = \int_a^b w(x) f(x) g(x) dx, \tag{9.4}$$

where w is a **weight function**, defined, positive, continuous and integrable on the open interval (a, b) . The norm $\| \cdot \|$, induced by this inner product and given by

$$\|f\| = \left(\int_a^b w(x) f^2(x) dx \right)^{1/2}, \tag{9.5}$$

is referred to as the 2-norm on $C[a, b]$ (see Example 8.2). For the sake of simplicity, we have chosen not to distinguish in terms of our notation between the 2-norm on $C[a, b]$ defined above and the 2-norm for vectors introduced in Chapter 2; it will always be clear from the context which of the two is intended.

Clearly, it is not necessary to demand the continuity of the function f on the closed interval $[a, b]$ to ensure that $\|f\|$ is finite. For example, $f: x \mapsto \text{sgn } x - (a + b)$, $x \in [a, b]$, has finite 2-norm, despite the fact that it has a jump discontinuity at $x = -(a + b)$.

Motivated by this observation, and the desire to develop a theory of approximation in the 2-norm whose range of applicability extends beyond the linear space of continuous functions on a bounded closed interval, we denote by $L^2(a, b)$ the set of all real-valued functions f

defined on (a, b) such that $w(x)f(x)$ is integrable on (a, b) ; the set $L_2(a, b)$ is equipped with the inner product (9.4) and the induced 2-norm (9.5). Obviously, $C[a, b]$ is a proper subset of $L_2(a, b)$.

In this broader context, $L_2(a, b)$ is frequently referred to as the L_2 -norm; for the sake of simplicity we shall continue to call it the 2-norm. As before, w is assumed to be a real-valued function, defined, positive, continuous and integrable on the open interval (a, b) . When $w(x) \equiv 1$ on (a, b) , we shall write $L_2(a, b)$ instead of $L_2(a, b)$.

We are now ready to consider best approximation in the 2-norm.

9.3 Best approximation in the 2-norm

The problem of best approximation in the 2-norm can be formulated as follows:

(B) Given that $f \in L_2(a, b)$, find p_n & ϵ_n such that

$$\|f - p_n\|_2 = \inf_{p \in P_n} \|f - p\|_2;$$

such p_n is called a **polynomial of best approximation of degree n to the function f in the 2-norm on (a, b)** .

The existence and uniqueness of p_n will be shown in Theorem 9.2. However, we shall first consider some simple examples.

Example 9.4 Suppose that $\epsilon > 0$ and let $f(x) = 1 - e^{-x}$ with x in $[0, 1]$. For $\epsilon = 10^{-6}$, the function f is depicted in Figure 9.1. We shall construct the polynomial of best approximation of degree 0 in the 2-norm, with weight function $w(x) \equiv 1$, for f on $(0, 1)$, and compare it with the minimax polynomial of degree 0 for f on $[0, 1]$.

The best approximation to f by a polynomial of degree 0 in the 2-norm on the interval $(0, 1)$, with weight function $w(x) \equiv 1$, is determined by minimising $\|f - c\|_2$ over all $c \in \mathbb{R}$; equivalently, we need to minimise

$$\|f(x) - c\|_2^2 = \int_0^1 (f(x) - c)^2 dx = \int_0^1 f(x)^2 dx - 2c \int_0^1 f(x) dx + c^2$$

$$\begin{aligned} &= \int_0^1 (1 - e^{-x})^2 dx - 2c \int_0^1 (1 - e^{-x}) dx + c^2 \\ &= \left[x - 2e^{-x} - x^2 \right]_0^1 - 2c \left[x + e^{-x} \right]_0^1 + c^2 \\ &= \left(1 - 2e^{-1} - 1 \right) - 2c \left(1 + e^{-1} \right) + c^2 \\ &= -2e^{-1} - 2c \left(1 + e^{-1} \right) + c^2 \end{aligned}$$

which tends to $\frac{1}{2}$ as $n \rightarrow \infty$. These examples indicate that the polynomial of best approximation from \mathcal{P}_n &

where

$$M_{kj} = \int_a^b x^k x^j dx = \frac{1}{k+j+1},$$

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

Equivalently, recalling that the inner product associated with the 2-norm (in the case of $w(x) = 1$) is defined by

$$(g, h) = \int_a^b g(x)h(x)dx,$$

M and b can be written as

$$M_{kj} = (x^k, x^j), \quad b_j = (f, x^j). \quad (9.7)$$

By solving the system of linear equations (9.6) for c_0, \dots, c_n , we obtain the coefficients of the polynomial of best approximation of degree n to the function f in the 2-norm on the interval $(0, 1)$. We can proceed in the same manner on any interval (a, b) with any positive, continuous and integrable weight function w defined on (a, b) .

This approach is straightforward for small values of n , but soon becomes impractical as n increases. The source of the computational difficulties is the fact that the matrix M is the Hilbert matrix, discussed in Section 2.8. The Hilbert matrix is well known to be ill-conditioned for large n , so any solution to (9.6), computed with a fixed number of decimal digits, loses all accuracy due to accumulation of rounding errors. Fortunately, an alternative method is available, and is discussed in the next section.

9.4 Orthogonal polynomials

In the previous section we described a method for constructing the polynomial of best approximation p_n to a function f in the 2-norm; it was based on seeking p_n as a linear combination of the polynomials x^j , $j = 0, \dots, n$, which form a basis for the linear space \mathcal{P}_n . The approach was not entirely satisfactory because it gave rise to a system of linear equations with a full matrix that was difficult to invert. The central idea of the alternative approach that will be described in this section is to expand p_n in terms of a different basis, chosen so that the resulting system of linear equations has a diagonal matrix; solving this

linear system is then a trivial exercise. Of course, the nontrivial ingredient of this alternative approach is to find a suitable basis for \mathcal{P}_n which achieves the objective that the matrix of the linear system is diagonal. The expression for M_{jk} in (9.7) gives us a clue how to proceed.

Suppose that $\phi_j, j = 0, \dots, n$, form a basis for \mathcal{P}_n , $n \geq 0$; let us seek the polynomial of best approximation as

$$p(x) = \sum_{j=0}^n c_j \phi_j(x),$$

where c_0, \dots, c_n are real numbers to be determined. By the same process as in the previous section, we arrive at a system of linear equations of the form (9.6):

$$M_{jk} c_k = f_j, \quad j = 0, \dots, n,$$

where now

$$M_{jk} = \int_a^b \phi_j(x) \phi_k(x) w(x) dx \quad \text{and} \quad f_j = \int_a^b \phi_j(x) f(x) w(x) dx,$$

with the inner product (g, h) defined by

$$(g, h) = \int_a^b w(x)g(x)h(x) dx,$$

and the weight function w assumed to be positive, continuous and integrable on the interval (a, b) .

Thus, $M = (M_{jk})$ will be a diagonal matrix provided that the basis functions $\phi_j, j = 0, \dots, n$, for the linear space \mathcal{P}_n are chosen so that $(\phi_j, \phi_k) = 0$, for $j \neq k$; in other words, $\{\phi_j\}$ is required to be orthogonal to $\{\phi_k\}$ for $j \neq k$, in the sense of Definition 9.2. This observation motivates the following definition.

Definition 9.4 Given a weight function w , defined, positive, continuous and integrable on the interval (a, b) , we say that the sequence of polynomials $\phi_j, j = 0, 1, \dots$, is a **system of orthogonal polynomials** on the interval (a, b) with respect to w , if each ϕ_j is of exact degree j , and if

$$\int_a^b w(x) \phi_j(x) \phi_k(x) dx = 0 \quad \text{for all } k \neq j, \\ = \int_a^b w(x) \phi_j^2(x) dx \quad \text{when } k = j.$$

Next, we show that a system of orthogonal polynomials exists on any interval (a, b) and for any weight function w which satisfies the conditions in Definition 9.4. We proceed inductively.

Let $p_j(x)$ be a polynomial of degree j , and suppose that p_0, \dots, p_{j-1} has already been constructed for $j = 0, \dots, n$, with $p_0(x) = 1$. Then,

$$\int_a^b p_j(x) p_k(x) w(x) dx = 0, \quad k = 0, \dots, j-1.$$

Let us now define the polynomial

$$q_j(x) = p_j(x) - a_0 p_0(x) - \dots - a_{j-1} p_{j-1}(x),$$

where

$$a_k = \frac{\int_a^b p_k(x) p_j(x) w(x) dx}{\int_a^b p_k(x)^2 w(x) dx}, \quad k = 0, \dots, j-1.$$

It then follows that

$$\begin{aligned} \int_a^b p_j(x) q_j(x) w(x) dx &= \int_a^b p_j(x) p_j(x) w(x) dx - \sum_{k=0}^{j-1} a_k \int_a^b p_j(x) p_k(x) w(x) dx \\ &= \int_a^b p_j(x)^2 w(x) dx - \sum_{k=0}^{j-1} a_k \int_a^b p_k(x)^2 w(x) dx \\ &= 0 \quad \text{for } 0 \leq k < j, \end{aligned}$$

where we have used the orthogonality of the sequence p_0, \dots, p_{j-1} , $j = 0, \dots, n$. Thus, with this choice of the numbers a_k we have ensured that q_j is orthogonal to all the previous members of the sequence, and p_{j+1} can now be defined as any nonzero-constant multiple of q_j . This procedure for constructing a system of orthogonal polynomials is usually referred to as **Gram–Schmidt orthogonalisation**.

Example 9.5 We shall construct a system of orthogonal polynomials p_0, p_1, \dots on the interval $(0, 1)$ with respect to the weight function $w(x) = 1$.

We put $p_0(x) = 1$, and we seek p_1 in the form

$$p_1(x) = x + c$$

such that $\int_0^1 p_1(x) p_0(x) dx = 0$; that is,

$$\int_0^1 (x + c) dx = 0.$$

¹ KX F B & K #. > / " C @ % # @ # ? \$
 / "DL B#* K 6 #. &? / B
 D C ? / , # @ > @ F 6D4

Hence,

$$c = \frac{x, \quad !}{, \quad !} = -$$

and therefore,

$$(x) = x \quad - \quad (x) = x \quad - .$$

By construction, $, \quad ! = , \quad ! = 0$.

We now seek \quad in the form

$$(x) = x \quad (d \quad (x) + d \quad (x))$$

such that $, \quad ! = 0$ and $, \quad ! = 0$. Thus,

$$\begin{aligned} x, \quad ! \quad d, \quad ! \quad d, \quad ! &= 0, \\ x, \quad ! \quad d, \quad ! \quad d, \quad ! &= 0. \end{aligned}$$

As $, \quad ! = 0$ and $, \quad ! = 0$, we have that

$$\begin{aligned} d &= \frac{x, \quad !}{, \quad !} = 1, \\ d &= \frac{x, \quad !}{, \quad !} = - , \end{aligned}$$

and therefore

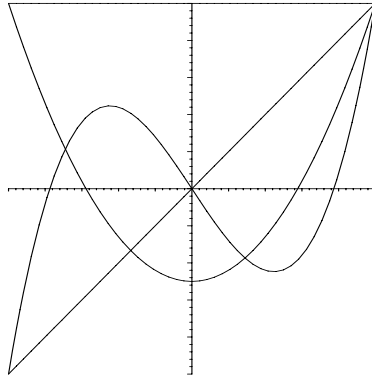
$$(x) = x \quad x + - . \tag{9.8}$$

Clearly, $, \quad ! = 0$ for $j = k, j, k \quad 0, 1, 2$, and \quad is of exact degree $j, j = 0, 1, 2$. Thus we have found the required system $, \quad , \quad$ of orthogonal polynomials on the interval $(0, 1)$ with respect to the given weight function w .

By continuing this procedure, we can construct a system of orthogonal polynomials $, \quad , \dots, \quad$, with respect to the weight function $w(x) \quad 1$ on the interval $(0, 1)$, for any $n \quad 1$. For example, when $n = 3$, we shall find $, \quad , \quad , \quad$, with $, \quad , \quad$, as above, and

$$(x) = x \quad -x \quad + -x \quad - .$$

Having generated a system of orthogonal polynomials on the interval $(0, 1)$ with respect to the weight function $w(x) \quad 1$, by performing the linear mapping $x \quad (b \quad a)x + a$ we may obtain a system of orthogonal polynomials on any open interval (a, b) with respect to the weight function $w(x) \quad 1$. For example, when $(a, b) = (-1, 1)$, the mapping $x \quad 2x \quad 1$ leads to the system of Legendre polynomials on $(-1, 1)$.



\$ 5 3 B * +

Example 9.6 (Legendre polynomials) We wish to construct a system of orthogonal polynomials on $(a, b) = (-1, 1)$ with respect to the weight function $w(x) = 1$.

On replacing x by

$$\frac{x - a}{b - a} = -(x + 1), \quad x \in (a, b) = (-1, 1),$$

in (x) , (x) , (x) , (x) from Example 9.5, we obtain, on normalising each of these polynomials so that its value at $x = 1$ is equal to 1, the polynomials P_0, P_1, P_2, P_3 , defined by

$$\begin{aligned} P_0(x) &= 1, \\ P_1(x) &= x, \\ P_2(x) &= \frac{1}{2}(3x^2 - 1), \\ P_3(x) &= \frac{1}{2}(5x^3 - 3x). \end{aligned}$$

These are the first four elements of the system of Legendre polynomials, orthogonal on the interval $(-1, 1)$ with respect to the weight function $w(x) = 1$. They are depicted in Figure 9.2. An alternative normalisation would have been to divide each P_n by $\sqrt{\int_{-1}^1 P_n^2(x) dx}$ so as to ensure that the 2-norm of the resulting scaled polynomial is equal to 1.

Example 9.7 The Chebyshev polynomials $T_n : x \mapsto \cos(n \cos^{-1} x)$, $n = 0, 1, \dots$, introduced in Section 8.4, form an orthogonal system on the interval $(-1, 1)$ with respect to the positive, continuous and integrable weight function $w(x) = (1 - x^2)^{-1/2}$.

The proof of this is simple: let (\cdot, \cdot) denote the inner product in $L_2(-1, 1)$

Such a system of polynomials is said to be **orthonormal**. The polynomials $\{p_j(x)\}_{j=0, \dots, n}$, $j = 0, \dots, n$, are linearly independent and form a basis for the linear space $C[a, b]$; therefore, each element $q \in C[a, b]$ can be expressed as a suitable linear combination,

$$q(x) = c_0 p_0(x) + c_1 p_1(x) + \dots + c_n p_n(x).$$

We wish to choose $\{c_j\}_{j=0, \dots, n}$, so as to ensure that the corresponding polynomial q minimises $\|f - q\|_2$ over all $q \in \mathcal{P}_n$. Let us, therefore, consider the function $E: (c_0, \dots, c_n) \mapsto E(c_0, \dots, c_n)$ defined by $E(c_0, \dots, c_n) = \|f - q\|_2^2$, where $q(x) = c_0 p_0(x) + c_1 p_1(x) + \dots + c_n p_n(x)$. Then,

$$\begin{aligned} E(c_0, \dots, c_n) &= \int_a^b (f(x) - q(x))^2 dx \\ &= \int_a^b (f(x) - c_0 p_0(x) - c_1 p_1(x) - \dots - c_n p_n(x))^2 dx \\ &= \int_a^b (f(x) - c_0 p_0(x) - c_1 p_1(x) - \dots - c_n p_n(x))^2 dx \\ &= \int_a^b (f(x) - c_0 p_0(x) - c_1 p_1(x) - \dots - c_n p_n(x))^2 dx \\ &= \int_a^b (f(x) - c_0 p_0(x) - c_1 p_1(x) - \dots - c_n p_n(x))^2 dx \end{aligned}$$

The function $(c_0, \dots, c_n) \mapsto E(c_0, \dots, c_n)$ achieves its minimum value at (c_0, \dots, c_n) , where

$$c_j = \frac{\int_a^b f(x) p_j(x) dx}{\int_a^b p_j(x)^2 dx}, \quad j = 0, \dots, n.$$

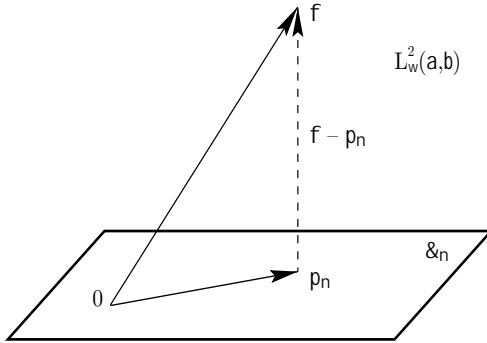
Hence $p \in \mathcal{P}_n$ defined by

$$p(x) = c_0 p_0(x) + c_1 p_1(x) + \dots + c_n p_n(x)$$

is the unique polynomial of best approximation of degree n to the function $f \in L^2(a, b)$ in the 2-norm on the interval (a, b) . \square

Remark 9.1 As $E(c_0, \dots, c_n) = \|f - p\|_2^2 = 0$, it follows from the proof of Theorem 9.2 that if $f \in L^2(a, b)$, and $\{p_j(x)\}_{j=0, \dots, n}$ is an orthonormal system of polynomials in $L^2(a, b)$, then

$$c_j = \int_a^b f(x) p_j(x) dx$$



$\int_a^b (f - p_n)^2 dx = \int_a^b (f - p_n + p_n - p_n)^2 dx$
 $= \int_a^b (f - p_n)^2 dx + \int_a^b (p_n - p_n)^2 dx + 2 \int_a^b (f - p_n)(p_n - p_n) dx$
 $= \int_a^b (f - p_n)^2 dx + 0 + 0$
 $= \int_a^b (f - p_n)^2 dx$

for each $n \geq 0$. This result is known as **Bessel's inequality**.

The next theorem, in conjunction with the use of orthogonal polynomials, will be our key tool for constructing the polynomial of best approximation in the 2-norm.

Theorem 9.3 A polynomial $p_n \in \mathcal{P}_n$ is the polynomial of best approximation of degree n to a function $f \in L^2_w(a,b)$ in the 2-norm if, and only if, the difference $f - p_n$ is orthogonal to every element of \mathcal{P}_n , i.e.,

$$\int_a^b (f - p_n)q \, dx = 0 \quad \forall q \in \mathcal{P}_n. \tag{9.9}$$

A geometrical illustration of the property (9.9) is given in Figure 9.3.

Proof of theorem Suppose that (9.9) holds. Then,

$$\int_a^b (f - p_n)q \, dx = 0 \quad \forall q \in \mathcal{P}_n,$$

given that $p_n \in \mathcal{P}_n$ for each $q \in \mathcal{P}_n$. Therefore,

$$\begin{aligned} \int_a^b (f - p_n)^2 \, dx &= \int_a^b (f - p_n)(f - p_n) \, dx \\ &= \int_a^b (f - p_n)(f - p_n + p_n - p_n) \, dx \\ &= \int_a^b (f - p_n)(f - p_n) \, dx + \int_a^b (f - p_n)(p_n - p_n) \, dx \\ &= \int_a^b (f - p_n)(f - p_n) \, dx + 0 \\ &= \int_a^b (f - p_n)^2 \, dx \end{aligned}$$

Hence, by the Cauchy–Schwarz inequality (9.2),

$$\|f - p\|_2 \leq \|f - q\|_2 \quad \forall q \in \mathcal{P}_n.$$

This implies that

$$\|f - p\|_2 \leq \|f - q\|_2 \quad \forall q \in \mathcal{P}_n.$$

On choosing $q = p$ on the right-hand side, equality will hold and therefore

$$\|f - p\|_2 = \min_{\mathcal{P}_n} \|f - q\|_2.$$

Conversely, suppose that p is the polynomial of best approximation to $f \in L_2(a, b)$. We have seen in the proof of Theorem 9.2 that p can be written in terms of the orthonormal polynomials ϕ_k , $k = 0, \dots, n$, as

$$p(x) = \sum_{k=0}^n c_k \phi_k(x),$$

where

$$c_k = \langle f, \phi_k \rangle, \quad k = 0, \dots, n. \tag{9.10}$$

On recalling that $\langle \phi_j, \phi_k \rangle = \delta_{j,k}$, $j, k = 0, \dots, n$, where $\delta_{j,k}$ is the Kronecker delta, we deduce from (9.10) that

$$\begin{aligned} \|f - p\|_2^2 &= \langle f - p, f - p \rangle = \langle f, f \rangle - 2 \sum_{k=0}^n c_k \langle f, \phi_k \rangle + \sum_{k=0}^n c_k^2 \\ &= \langle f, f \rangle - 2 \sum_{k=0}^n c_k^2 + \sum_{k=0}^n c_k^2 \\ &= \langle f, f \rangle - \sum_{k=0}^n c_k^2 = 0, \quad j = 0, \dots, n. \end{aligned} \tag{9.11}$$

Since $\mathcal{P}_n = \text{span}\{\phi_0, \dots, \phi_n\}$, it follows from (9.11) that $\|f - p\|_2 = 0$ for all $q \in \mathcal{P}_n$, as required. \square

An equivalent, but slightly more explicit, form of writing (9.9) is

$$\int_a^b w(x)(f(x) - p(x))q(x) dx = 0 \quad \forall q \in \mathcal{P}_n.$$

Theorem 9.2 provides a simple method for determining the polynomial of best approximation $p \in \mathcal{P}_n$ to a function $f \in L_2(a, b)$ in the 2-norm. First, proceeding as described in the discussion following Definition 9.4, we construct the system of orthogonal polynomials ϕ_j , $j = 0, \dots, n$, on the interval (a, b) with respect to the weight function w , if this system

is not already known. We normalise the polynomials $p_j(x)$, $j = 0, \dots, n$, by setting

$$p_j(x) = \frac{f_j(x)}{\|f_j\|}, \quad j = 0, \dots, n,$$

to obtain the system of orthonormal polynomials $p_j(x)$, $j = 0, \dots, n$, on (a, b) . We then evaluate the coefficients $c_j = \int_a^b f(x) p_j(x) dx$, $j = 0, \dots, n$, and form $p(x) = \sum_{j=0}^n c_j p_j(x)$.

We may avoid the necessity of determining the normalised polynomials by writing

$$\begin{aligned} p(x) &= \sum_{j=0}^n c_j p_j(x) \\ &= \sum_{j=0}^n \frac{c_j}{\|f_j\|} f_j(x) \\ &= \sum_{j=0}^n \frac{f_j(x)}{\|f_j\|} c_j, \end{aligned} \tag{9.12}$$

where

$$c_j = \frac{\int_a^b f(x) f_j(x) dx}{\int_a^b f_j^2(x) dx}, \quad j = 0, \dots, n. \tag{9.13}$$

Thus, as indicated at the beginning of the section, with this approach to the construction of the polynomial of best approximation in the 2-norm, we obtain the coefficients explicitly and there is no need to solve a system of linear equations with a full matrix.

Example 9.8 We shall construct the polynomial of best approximation of degree 2 in the 2-norm to the function $f: x \mapsto e^{-x}$ over $(0, 1)$ with weight function $w(x) = 1$.

We already know a system of orthogonal polynomials p_0, p_1, p_2 on this interval from Example 9.5; thus, we seek p in the form

$$p(x) = c_0 p_0(x) + c_1 p_1(x) + c_2 p_2(x), \tag{9.14}$$

where, according to (9.13),

$$c_j = \frac{\int_0^1 e^{-x} p_j(x) dx}{\int_0^1 p_j^2(x) dx}, \quad j = 0, 1, 2.$$

Recalling from Example 9.5 that

$$p_0(x) = 1, \quad p_1(x) = x - \frac{1}{2}, \quad p_2(x) = x^2 - x + \frac{1}{6},$$

we then have that

$$\begin{aligned}
 &= \frac{e}{1} = e - 1, \\
 &= \frac{3/2 \cdot e/2}{1/12} = 18 - 6e, \\
 &= \frac{7e/6 - 19/6}{1/180} = 210e - 570.
 \end{aligned}
 \tag{9.15}$$

Substituting the values of a_0 , a_1 , and a_2 into (9.14), we conclude that the polynomial of best approximation of degree 2 for the function $f: x \mapsto e^x$ in the 2-norm is

$$p_2(x) = (210e - 570)x^2 + (588 - 216e)x + (39e - 105).$$

The approximation error is

$$\|f - p_2\|_2 = 0.005431,$$

to six decimal digits.

We conclude this section by giving a property of orthogonal polynomials that will be required in the next chapter.

Theorem 9.4 *Suppose that $\{p_j\}_{j=0}^{\infty}$, $j = 0, 1, \dots$, is a system of orthogonal polynomials on the interval (a, b) with respect to the positive, continuous and integrable weight function w on (a, b) . It is understood that p_j is a polynomial of exact degree j . Then, for $j \geq 1$, the zeros of the polynomial p_j are real and distinct, and lie in the interval (a, b) .*

Proof Suppose that x_i , $i = 1, \dots, k$, are the points in the open interval (a, b) at which $p_j(x)$ changes sign. Let us note that $k \geq 1$, because for $j \geq 1$, by orthogonality of $p_j(x)$ to $p_0(x) = 1$, we have that

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

Thus, the integrand, being a continuous function that is not identically zero on (a, b) , must change sign on (a, b) ; however, w is positive on (a, b) , so p_j must change sign at least once on (a, b) . Therefore $k \geq 1$.

Let us define

$$q(x) = (x - x_1) \dots (x - x_k). \tag{9.16}$$

Now the function $p_j(x) / q(x)$ does not change sign in the interval (a, b) , since at each point where $p_j(x)$ changes sign $q(x)$ changes sign also. Hence,

$$\int_a^b w(x) \frac{p_j(x)}{q(x)} dx = 0.$$

However, ϕ_j is orthogonal to every polynomial of lower degree with respect to the weight function w , so the degree of the polynomial must be at least j ; thus, $k \geq j$. On the other hand, k cannot be greater than j , since a polynomial of exact degree j cannot change sign more than j times. Therefore $k = j$; *i.e.*, the points (a, b) , $i = 1, \dots, j$, are the zeros (and all the zeros) of $\phi_j(x)$. \square

9.5 Comparisons

We can show that the polynomial of best approximation in the 2-norm for a function $f \in C[a, b]$ is also a near-best approximation in the ∞ -norm for f on $[a, b]$ in the sense defined in Section 8.5.

Theorem 9.5 *Let $n \geq 0$ and assume that f is defined and continuous on the interval $[a, b]$, and $f \in C[a, b]$. Let p_n be the polynomial of best approximation of degree n to f in the 2-norm on $[a, b]$, where the weight function w is positive, continuous and integrable on (a, b) . Then, the difference $f - p_n$ changes sign at no less than $n + 1$ distinct points in the interval (a, b) .*

Proof The proof is very similar to that of Theorem 9.4; we shall give an outline and leave the details as an exercise.

As $\int_a^b (f - p_n)^2 w(x) dx = 0$, *i.e.*,

$$\int_a^b w(x)(f(x) - p_n(x)) dx = 0,$$

and $w(x) > 0$ for all $x \in (a, b)$, it follows that $f - p_n$ changes sign in (a, b) . Let ξ_j , $j = 1, \dots, k$, denote distinct points in (a, b) where $f - p_n$ changes sign. We shall prove that $k \geq n + 1$.

Define the polynomial $\phi_n(x)$ as in (9.16); then, $w(x)[f(x) - p_n(x)] \phi_n(x)$ does not change sign in (a, b) , and so its integral over (a, b) is not zero. Therefore, $\int_a^b (f - p_n) \phi_n w(x) dx \neq 0$. On the other hand, according to Theorem 9.3, $f - p_n$ is orthogonal to every polynomial of degree n or less. Hence the degree of $\phi_n(x)$ must be greater than n , and so $k \geq n + 1$. \square

We return to the example illustrated by Figure 8.5, and consider the difference $f - p$ for the function $f: x \rightarrow e^{-x}$ on the interval $(0, 1)$. Figure 9.4 shows this difference for two polynomial approximations of degree 4: the minimax approximation of Section 8.5 and the best approximation in the 2-norm with weight function $w(x) = 1$. It is clear that the

which gives greater weight near the ends of the interval, it seems likely that the extrema of the error might be more nearly equal. This can be achieved by using the weight function $w(x) = [x(1-x)]^{-1/2}$, so that the orthogonal polynomials are the Chebyshev polynomials adapted to the interval $(0, 1)$. Figure 9.5 shows the corresponding difference $f - p$, and we now see that the two best approximations, in the ∞ -norm and the weighted 2-norm, are very close.

Polynomials of best approximation in the 2-norm have a special property which is often useful. Suppose that we have constructed the best polynomial approximation, p_n , of degree n , in the 2-norm, but that p_n does not achieve the required accuracy. To construct the best polynomial approximation of degree $n + 1$ all we need is to calculate

$$= \frac{\int_0^1 (f - p_n)^2 w(x) dx}{\int_0^1 w(x) dx}$$

and then let $p_{n+1}(x) = p_n(x) + c_{n+1} T_{n+1}(x)$. By noting that

$$\int_0^1 (f - p_n - c_{n+1} T_{n+1})^2 w(x) dx = 0, \quad j = 0, 1, \dots, n+1,$$

it follows that p_{n+1} is best least squares approximation to f from \mathcal{P}_{n+1} . If we are constructing the minimax approximation of degree $n + 1$, or using Lagrange interpolation with equally spaced points, the work involved in constructing p_n is lost, and the construction of p_{n+1} must begin completely afresh.

9.6 Notes

We give some pointers to the vast literature on orthogonal polynomials. The following are classical sources on the subject.

1. G. G. Lorentz, "Orthogonal Polynomials", Pergamon Press, Oxford, New York, 1971.

\$ 1 "Orthogonal Polynomials", Memoirs of the American Mathematical Society, no. 213, American Mathematical Society, Providence, RI, 1979.

1 & 2 "Orthogonal Polynomials", Colloquium publications (American Mathematical Society), 23, American Mathematical Society, Providence, RI, 1959.

Tables of orthogonal polynomials are found in

+) (. #) & (Editors), 'Orthogonal polynomials', Ch. 22 in *Handbook of Mathematical Functions with Formulas, Graphs, and Mathematical Tables*, ninth printing, Dover, New York, pp. 771–802, 1972.

Computational aspects of the theory of orthogonal polynomials are discussed in the edited volume

" " ' (Editors), *Applications and Computation of Orthogonal Polynomials*, Conference at the Mathematical Research Institute, Oberwolfach, Germany, March 22–28, 1998, Birkhäuser, Basel, 1999.

A recent survey of the theory and application of orthogonal polynomials in numerical computations is contained in

" Orthogonal polynomials: applications and computation, *Acta Numerica* **5** (A. Iserles, ed.), Cambridge University Press, Cambridge, pp. 45–119, 1996.

Finally, we refer to the books of Powell and Cheney, cited in the Notes at the end of the previous chapter, concerning the application of orthogonal polynomials in the field of best least squares approximation.

Exercises

- 9.1 Construct orthogonal polynomials of degrees 0, 1 and 2 on the interval $(0, 1)$ with the weight function $w(x) = \ln x$.
- 9.2 Let the polynomials p_j , $j = 0, 1, \dots$, form an orthogonal system on the interval $(-1, 1)$ with respect to the weight function $w(x) = 1$. Show that the polynomials $p_j((2x - a - b)/(b - a))$, $j = 0, 1, \dots$, represent an orthogonal system for the interval (a, b) and the same weight function. Hence obtain the polynomials in Example 9.5 from the Legendre polynomials in Example 9.6.
- 9.3 Suppose that the polynomials p_j , $j = 0, 1, \dots$, form an orthogonal system on the interval $(0, 1)$ with respect to the weight function $w(x) = x^\alpha$, $\alpha > 0$. Find, in terms of p_j , a system of orthogonal polynomials for the interval $(0, b)$ and the same weight function.

9.4 Show, by induction or otherwise, that, for $0 \leq k < n$,

$$\frac{d}{dx} (1 - x^2)^n = (1 - x^2)^{n-1} q(x),$$

where q is a polynomial of degree k . Deduce that all the derivatives of the function $(1 - x^2)^n$ of order less than n vanish at $x = \pm 1$.

Define $P_k(x) = (d^k/dx^k) (1 - x^2)^n$, and show by repeated integration by parts that

$$\int_{-1}^1 P_k(x) P_j(x) dx = 0, \quad 0 \leq k < j.$$

Hence verify the expressions in Example 9.6 for the Legendre polynomials of degrees 0, 1, 2 and 3.

9.5 Show, by induction or otherwise, that, for $0 \leq k < j$,

$$\frac{d}{dx} (x^k e^{-x}) = x^{k-1} q(x) e^{-x},$$

where $q(x)$ is a polynomial of degree k .

The function $P_j(x)$ is defined for $j \geq 0$ by

$$P_j(x) = e^{-x} \frac{d^j}{dx^j} (x^j e^{-x}).$$

Show that, for each $j \geq 0$, $P_j(x)$ is a polynomial of degree j , and that these polynomials form an orthogonal system on the interval $(0, \infty)$ with respect to the weight function $w(x) = e^{-x}$. Write down the polynomials with $j = 0, 1, 2$ and 3.

9.6 Suppose that $P_j, j = 0, 1, \dots$, form a system of orthogonal polynomials with weight function $w(x)$ on the interval (a, b) . Show that, for some value of the constant C , $(x - C) P_j(x)$ is a polynomial of degree j , and hence that

$$(x - C) P_j(x) = P_j(x) + D_j P_{j-1}(x),$$

Use the orthogonality properties to show that $D_j = 0$ for $k < j - 1$, and deduce that the polynomials satisfy a recurrence relation of the form

$$(x - C) P_j(x) + D_{j+1} P_{j+1}(x) + E_j P_{j-1}(x) = 0, \quad j \geq 1.$$

9.7 In the notation of Exercise 6 suppose that the normalisation of the polynomials is so chosen that for each j the coefficient of x^j in $w(x)$ is positive. Show that $C_j > 0$ for all j . By considering

$$\int_a^b w(x) (x^j - C_j) (x^j - C_j) dx$$

show that

$$\int_a^b w(x) x^j - C_j (x^j) dx > 0,$$

and deduce that $C_j > 0$ for all j . Hence show that for all positive values of j the zeros of p_n and p_{n-1} interlace. (See the proof of Theorem 5.8.)

9.8 Using the weight function w on the interval (a, b) apply a similar argument to that for Theorem 8.6 to find the best polynomial approximation p_n of degree n in the 2-norm to the function x^c . Show that

$$x^c - p_n = \frac{\int_a^b w(x) x^c dx}{c} - p_n,$$

where c is the coefficient of x^c in $w(x)$.

Write down the best polynomial approximation of degree 2 to the function x^c in the 2-norm with $w(x) = 1$ on the interval $(-1, 1)$, and evaluate the 2-norm of the error.

9.9 Suppose that the weight w is an even function on the interval $(-a, a)$, and that a system of orthogonal polynomials p_j , $j = 0, \dots, n$, on the interval $(-a, a)$ is constructed by the Gram-Schmidt process. Show that, if j is even, then p_j is an even function, and that, if j is odd, then p_j is an odd function.

Now suppose that the best polynomial approximation of degree n in the 2-norm to the function f on the interval $(-a, a)$ is expressed in the form

$$p_n(x) = \sum_{j=0}^n c_j p_j(x) + \sum_{j=0}^n d_j p_j(x).$$

Show that if f is an even function, then all the odd coefficients d_j are zero, and that if f is an odd function, then all the even coefficients c_j are zero.

9.10 The function $H(x)$ is defined by $H(x) = 1$ if $x > 0$, and $H(-x) = H(x)$. Construct the best polynomial approximations of degrees 0, 1 and 2 in the 2-norm to this function over the interval $(-1, 1)$ with weight function $w(x) = 1$. (It may not

appear very useful to consider a polynomial approximation to a discontinuous function, but representations of such functions by Fourier series will be familiar to most readers. Note that the function H belongs to $L^2(-1, 1)$.

Numerical integration – II

10.1 Introduction

In Section 7.2 we described the Newton–Cotes family of formulae for numerical integration. These were constructed by replacing the integrand by its Lagrange interpolation polynomial with equally spaced interpolation points and integrating this exactly. Here, we consider another family of numerical integration rules, called Gauss quadrature formulae, which are based on replacing the integrand f by its Hermite interpolation polynomial and choosing the interpolation points x_i in such a way that, after integrating the Hermite polynomial, the derivative values $f'(x_i)$ do not enter the quadrature formula. It turns out that this can be achieved by requiring that the x_i are roots of a polynomial of a certain degree from a system of orthogonal polynomials.

10.2 Construction of Gauss quadrature rules

Suppose that the function f is defined on the closed interval $[a, b]$ and that it is continuous and differentiable on this interval. Suppose, further, that w is a weight function, defined, positive, continuous and integrable on (a, b) . We wish to construct quadrature formulae for the approximate evaluation of the integral

$$\int_a^b w(x)f(x)dx.$$

For a nonnegative integer n , let x_i , $i = 0, \dots, n$, be $n + 1$ points in the interval $[a, b]$; the precise location of these points will be determined later on. The Hermite interpolation polynomial of degree $2n + 1$ for the

function f is given by the expression (see Section 6.4)

$$p_n(x) = H(x)f(x) + K(x)f'(x), \quad (10.1)$$

where

$$\begin{aligned} H(x) &= [L(x)](1 - 2L(x)(x - x_0)), \\ K(x) &= [L(x)](x - x_0). \end{aligned} \quad (10.2)$$

Further, for $n \geq 1$, L_k is defined by

$$L_k(x) = \frac{x - x_{k+1}}{x - x_k}, \quad k = 0, 1, \dots, n;$$

if $n = 0$, we let $L_0(x) = 1$ and thereby $H(x) = 1$ and $K(x) = x - x_0$ for this value of n . Thus, we deduce from (10.1) that

$$\begin{aligned} \int_a^b w(x)f(x)dx &= \int_a^b w(x)p_n(x)dx \\ &= W \int_a^b f(x)dx + V \int_a^b f'(x)dx, \end{aligned} \quad (10.3)$$

where

$$W = \int_a^b w(x)H(x)dx, \quad V = \int_a^b w(x)K(x)dx.$$

There is an obvious advantage in choosing the points x_k in such a way that all the coefficients V_k are zero, for then the derivative values $f'(x_k)$ are not required. Recalling the form of the polynomial K and inserting it into the defining expression for V , we have

$$\begin{aligned} V &= \int_a^b w(x)[L(x)](x - x_0)dx \\ &= C \int_a^b w(x)(x - x_0)L(x)dx, \end{aligned} \quad (10.4)$$

where $L(x) = (x - x_1)\dots(x - x_n)$ and

$$C = \begin{cases} + & (x - x_0)^- & \text{if } n \geq 1, \\ 1 & & \text{if } n = 0. \end{cases}$$

Since $L(x)$ is of degree $n + 1$ while $L_k(x)$ is of degree n for each k , $0 \leq k \leq n$, each V_k will be zero if the polynomial $L_k(x)$ is orthogonal to every polynomial of lower degree with respect to the weight function

w . We can therefore construct the required quadrature formula (10.3) with $V = 0$, $k = 0, \dots, n$, by choosing the points x_k , $k = 0, \dots, n$, to be the zeros of the polynomial of degree $n + 1$ in a system of orthogonal polynomials over the interval (a, b) with respect to the weight function w ; we know from Theorem 9.4 that these zeros are real and distinct, and all lie in the open interval (a, b) .

Having chosen the location of the points x_k , we now consider W :

$$\begin{aligned} W &= \int_a^b w(x) H_n(x) dx \\ &= \int_a^b w(x) [L_{n+1}(x) - 2L_n(x)(x - x_n)] dx \\ &= \int_a^b w(x) [L_{n+1}(x) - 2L_n(x)] V_n(x) dx. \end{aligned} \quad (10.5)$$

Since $V_n = 0$, the second term in the last line vanishes and thus we obtain the following numerical integration formula, known as the **Gauss quadrature** rule:

$$\int_a^b w(x) f(x) dx \approx W \sum_{k=0}^n f(x_k), \quad (10.6)$$

where the **quadrature weights** are

$$W_k = \int_a^b w(x) [L_n(x)]^2 dx, \quad (10.7)$$

and the **quadrature points** x_k , $k = 0, \dots, n$, are chosen as the zeros of the polynomial of degree $n + 1$ from a system of orthogonal polynomials over the interval (a, b) with respect to the weight function w . Since this quadrature rule was obtained by exact integration of the Hermite interpolation polynomial of degree $2n + 1$ for f , it gives the exact result whenever f is a polynomial of degree $2n + 1$ or less.

Example 10.1 Consider the case $n = 1$, with the weight function $w(x) = 1$ over the interval $(0, 1)$.

The quadrature points x_0, x_1 are then the zeros of the polynomial constructed in Example 9.5 and given by (9.8),

$$(x) = x^2 - x + \frac{1}{4}, \quad (10.8)$$

and therefore

$$x_1 = -\frac{1}{\sqrt{3}}, \quad x_2 = \frac{1}{\sqrt{3}}.$$

Clearly, x_1 and x_2 belong to the open interval $(0, 1)$, in accordance with Theorem 9.4. The weights are obtained from (10.7):

$$\begin{aligned} W_1 &= \int_{-1}^1 \frac{x_1 - x}{x_1 - x} dx \\ &= 3 \int_{-1}^1 (x - 2x^2 + x^3) dx \\ &= 3(-\frac{1}{2} + \frac{2}{3} - \frac{1}{4}) \\ &= \frac{1}{3}, \end{aligned} \tag{10.9}$$

and $W_2 = \frac{1}{3}$ in the same way. We thus have the Gauss quadrature rule

$$\int_{-1}^1 f(x) dx \approx \frac{1}{3} f(-\frac{1}{\sqrt{3}}) + \frac{1}{3} f(\frac{1}{\sqrt{3}}), \tag{10.10}$$

which is exact whenever f is a polynomial of degree $2 - 1 + 1 = 3$ or less.

10.3 Direct construction

The calculation of the weights and the quadrature points in a Gauss quadrature rule requires little work when the system of orthogonal polynomials is already known. If this is not known, at the very least it is necessary to construct the polynomial from the system whose roots are the quadrature points; in that case a straightforward approach, which avoids this construction, may be easier.

Suppose, for example, that we wish to find the values of A_1, A_2, x_1 and x_2 such that the quadrature rule

$$\int_{-1}^1 f(x) dx \approx A_1 f(x_1) + A_2 f(x_2) \tag{10.11}$$

is exact for all $f \in \mathcal{P}_3$.

We have to determine four unknowns, A_1, A_2, x_1 and x_2 , so we need four equations; thus we take, in turn, $f(x) = 1$, $f(x) = x$, $f(x) = x^2$ and $f(x) = x^3$ and demand that the quadrature rule (10.11) is exact (that is, the integral of f is equal to the corresponding approximation obtained by inserting f into the right-hand side of (10.11)). Hence,

$$1 = A_1 + A_2, \tag{10.12}$$

$$- = A x + A x , \tag{10.13}$$

$$- = A x + A x , \tag{10.14}$$

$$- = A x + A x . \tag{10.15}$$

It remains to solve this system. To do so, we consider the quadratic polynomial defined by

$$(x) = (x - x_1)(x - x_2)$$

whose roots are the unknown quadrature points x_1 and x_2 . In expanded form, (x) can be written as

$$(x) = x^2 + px + q.$$

First we shall determine p and q ; then we shall find the roots x_1 and x_2 of (x) . We shall then insert the values of x_1 and x_2 into (10.13) and solve the linear system (10.12), (10.13) for A_1 and A_2 .

To find p and q , we multiply (10.12) by q , (10.13) by p and (10.14) by 1, and we add up the resulting equations to deduce that

$$\begin{aligned} - + -p + q &= A_1 (x_1 + px_1 + q) + A_2 (x_1 + px_1 + q) \\ &= A_1 (x_1) + A_2 (x_1) = A_1 \cdot 0 + A_2 \cdot 0 = 0. \end{aligned}$$

Therefore,

$$- + -p + q = 0. \tag{10.16}$$

Similarly, we multiply (10.13) by q , (10.14) by p and (10.15) by 1, and we add up the resulting equations to obtain

$$\begin{aligned} - + -p + -q &= A_1 x_2 (x_2 + px_2 + q) + A_2 x_2 (x_2 + px_2 + q) \\ &= A_1 x_2 (x_2) + A_2 x_2 (x_2) = A_1 \cdot 0 + A_2 \cdot 0 = 0. \end{aligned}$$

Thus,

$$- + -p + -q = 0. \tag{10.17}$$

From (10.16) and (10.17) we immediately find that $p = -1$ and $q = -$. Having determined p and q , we see that

$$(x) = x^2 - x - ,$$

in agreement with (10.8). We then find the roots of this quadratic polynomial to give x_1 and x_2 as before. With these values of x_1 and x_2 we deduce from (10.12) and (10.13) that

$$\begin{aligned} A_1 + A_2 &= 1, \\ A_1 (- + -) - A_2 (- -) &= 0, \end{aligned}$$

and therefore $A = A = -$. Thus, we conclude that the required quadrature rule is (10.10), as before.

It is easy to see that equations (10.16) and (10.17) express the condition that the polynomial $x^2 + px + q$ is orthogonal to the polynomials 1 and x respectively. This alternative approach has simply constructed a quadratic polynomial from a system of orthogonal polynomials by requiring that it is orthogonal to every polynomial of lower degree, instead of building up the whole system of orthogonal polynomials.

A straightforward calculation shows that, in general, the quadrature rule (10.10) is not exact for polynomials of degree higher than 3 (take $f(x) = x^3$, for example, to verify this).

Example 10.2 We shall apply the quadrature rule (10.10) to compute an approximation to the integral $I = \int_{-1}^1 e^x dx$.

Using (10.10) with $f(x) = \exp(x) = e^x$ yields

$$I \approx -\exp(-1) - \frac{1}{2} + -\exp(1) + \frac{1}{2} = \bar{e} \cosh \bar{1}.$$

On rounding to six decimal digits, $I \approx 1.717896$. The exact value of the integral is $I = e - 1 = 1.718282$, rounding to six decimal digits.

10.4 Error estimation for Gauss quadrature

The next theorem provides a bound on the error that has been committed by approximating the integral on the left-hand side of (10.6) by the quadrature rule on the right.

Theorem 10.1 Suppose that w is a weight function, defined, integrable, continuous and positive on (a, b) , and that f is defined and continuous on $[a, b]$; suppose further that f has a continuous derivative of order $2n + 2$ on $[a, b]$, $n \geq 0$. Then, there exists a number η in (a, b) such that

$$\int_a^b w(x)f(x)dx - W_n(f) = K f^{(2n+2)}(\eta), \tag{10.18}$$

and

$$K = \frac{1}{(2n+2)!} \int_a^b w(x)[f^{(2n+2)}(x)]^2 dx.$$

Consequently, the integration formula (10.6), (10.7) will give the exact result for every polynomial of degree $2n + 1$.

Proof Recalling the definition of the Hermite interpolation polynomial p for the function f and using Theorem 6.4, we have

$$\begin{aligned} \int_a^b w(x)f(x)dx - \int_a^b w(x)p(x)dx &= \int_a^b w(x)[f(x) - p(x)]dx \\ &= \int_a^b w(x) \frac{f^{(2n+2)}(\xi)}{(2n+2)!} [\omega(x)] dx. \end{aligned} \tag{10.19}$$

However, by the Integral Mean Value Theorem, Theorem A.6, the last term is equal to

$$\frac{f^{(2n+2)}(\xi)}{(2n+2)!} \int_a^b w(x)[\omega(x)] dx,$$

for some $\xi \in (a, b)$, and hence the desired error bound. □

Note that, by virtue of Theorem 10.1, the Gauss quadrature rule gives the exact value of the integral when f is a polynomial of degree $2n + 1$ or less, which is the highest possible degree that one can hope for with the $2n + 2$ free parameters consisting of the quadrature weights W_k , $k = 0, \dots, n$, and the quadrature points x_k , $k = 0, \dots, n$.

A different approach leads to a proof of convergence of the Gauss formulae $I_n(f)$, defined in (10.6), (10.7), as $n \rightarrow \infty$.

Theorem 10.2 *Suppose that the weight function w is defined, positive, continuous and integrable on the open interval (a, b) . Suppose also that the function f is continuous on the closed interval $[a, b]$. Then,*

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

Proof If we choose any positive real number ϵ then, since f is continuous on $[a, b]$, the Weierstrass Theorem (Theorem 8.1) shows that there is a polynomial p such that

$$|f(x) - p(x)| < \epsilon \quad \text{for all } x \in [a, b]. \tag{10.20}$$

Let N be the degree of this polynomial, and write p as p_N .

Thus we deduce that

$$\begin{aligned} \int_a^b w(x)f(x)dx - I_n(f) &= \int_a^b w(x)[f(x) - p_N(x)]dx \\ &+ \int_a^b w(x)p_N(x)dx - I_n(p_N) \\ &+ I_n(p_N) - I_n(f). \end{aligned} \tag{10.21}$$

Consider the first term on the right of this equality; it follows from (10.20) that

$$\int_a^b w(x)[f(x) - p(x)]dx = W_1,$$

where

$$W = \int_a^b w(x)dx.$$

For the last term on the right of (10.21),

$$\begin{aligned} \int_a^b (f(x) - p(x))w(x)dx &= W[f(x) - p(x)] \\ &= \int_a^b w(x)dx \\ &= W, \end{aligned} \tag{10.22}$$

where we have used the fact that all the quadrature weights W_i are positive (see (10.7)), and that a Gauss quadrature rule integrates a constant function exactly. Now for the middle term in (10.21), if we define N to be the integer part of $-N$, we see that when $n \geq N$ the quadrature formula is exact for all polynomials of degree $2N + 1$ or less, and hence for the polynomial p (given that $N \leq 2N + 1 \leq 2n + 1$). Therefore,

$$\int_a^b w(x)p(x)dx - \sum_{i=1}^n w(x_i)p(x_i) = 0 \quad \text{if } n \geq N.$$

Putting these three terms together, we see that

$$\int_a^b w(x)f(x)dx - \sum_{i=1}^n w(x_i)f(x_i) = W + 0 + W \quad \text{if } n \geq N.$$

Finally, given any positive number ϵ , we define $\delta = \epsilon/(2W)$ and find the corresponding value of $N = N(\delta)$ to deduce that

$$\left| \int_a^b w(x)f(x)dx - \sum_{i=1}^n w(x_i)f(x_i) \right| < \epsilon \quad \text{if } n \geq N,$$

which is what we were required to prove. \square

The interest of this theorem is mainly theoretical, as it gives no indication of how rapidly the error tends to zero. However, it does show the importance of the fact that the weights W are positive. Much of the above proof would apply with little change to the Newton–Cotes formulae of Section 7.2. We saw there that for the formulae of order 1 and 2, the trapezium rule and Simpson’s rule, the weights are positive. However, when $n > 8$ some of the weights in the Newton–Cotes formula of order n become negative. In this case we have $W = (b - a)$, but we find that W as n , so the proof breaks down. Stronger conditions must be imposed on the function f to ensure that the Newton–Cotes formula converges to the required integral. (See the example in Section 7.4.)

10.5 Composite Gauss formulae

It is often useful to define composite Gauss formulae, just as we did for the trapezium rule and Simpson’s rule in Section 7.5. Let us suppose, for the sake of simplicity, that $w(x) = 1$. We divide the range $[a, b]$ into m subintervals $[x_{j-1}, x_j]$, $j = 1, 2, \dots, m$, $m \geq 2$, each of width $h = (b - a)/m$, and write

$$\int_a^b f(x) dx = \sum_{j=1}^m \int_{x_{j-1}}^{x_j} f(x) dx,$$

where

$$x_j = a + jh, \quad j = 0, 1, \dots, m.$$

We then map each of the subintervals $[x_{j-1}, x_j]$, $j = 1, 2, \dots, m$, onto the reference interval $[-1, 1]$ by the change of variable

$$x = -(x_{j-1} + x_j) + ht, \quad t \in [-1, 1],$$

giving

$$\int_{x_{j-1}}^{x_j} f(x) dx = -h \int_{-1}^1 g(t) dt = -h \int_{-1}^1 f(-(x_{j-1} + x_j) + ht) dt,$$

where

$$g(t) = f(-(x_{j-1} + x_j) + ht) \quad \text{and} \quad I_j = \int_{x_{j-1}}^{x_j} f(x) dx.$$

The composite Gauss quadrature rule is then obtained by applying

the same Gauss formula to each of the integrals I_k . This gives

$$\begin{aligned} \int_a^b f(x) dx &= \sum_{k=0}^n W_k f(x_k) \\ &= \sum_{k=0}^n W_k f\left(\frac{b-a}{2} + x_k\right) + \text{error}, \end{aligned} \quad (10.23)$$

where x_k are the quadrature points in $(-1, 1)$ and W_k are the associated weights for $k = 0, \dots, n$ with $\sum_{k=0}^n W_k = 2$.

An expression for the error of this composite formula is obtained, as in Section 7.5, by adding the expressions (10.18) for the errors in the integrals I_k . The result is

$$\text{error} = C \frac{(b-a)^{n+2}}{2^{n+1} m^{n+1} (2n+2)!} f^{(n+2)}(\xi) \quad (10.24)$$

where $\xi \in (a, b)$ and

$$C = \frac{1}{2^{n+1} m^{n+1}} \int_{-1}^1 |t|^{n+2} dt.$$

Definition 10.1 *The composite midpoint rule is the composite Gauss formula with $w(x) = 1$ and $n = 0$ defined by*

$$\int_a^b f(x) dx \approx h \sum_{j=0}^n f(a + (j + \frac{1}{2})h). \quad (10.25)$$

This follows from the fact that when $n = 0$ there is one quadrature point $x_0 = 0$ in $(-1, 1)$, which is at the midpoint of the interval, and the corresponding quadrature weight W_0 is equal to the length of the interval $(-1, 1)$, i.e., $W_0 = 2$. It follows from (10.24) with $n = 0$ and

$$C = \frac{1}{2^2 m^2} \int_{-1}^1 |t|^2 dt = \frac{1}{6}$$

that the error in the composite midpoint rule is

$$\text{error} = \frac{(b-a)^3}{24m^2} f''(\xi),$$

where $\xi \in (a, b)$, provided that the function f has a continuous second derivative on $[a, b]$.

10.6 Radau and Lobatto quadrature

We have now discussed two types of quadrature formulae, which have the same form, $\sum W_k f(x_k)$. In the Newton–Cotes formulae the (equally spaced) quadrature points x_k are given, and we were able to find the weights W_k so that the result was exact for polynomials of degree n . By allowing the quadrature points as well as the weights to be freely chosen, we constructed Gauss quadrature formulae which were exact for polynomials of degree $2n + 1$. There are also many possible formulae of mixed type, where some, but not all, of the quadrature points are given, and the rest can be freely chosen. We might expect that each quadrature point which is fixed will reduce the degree of polynomial for which such a formula is exact by 1, from the maximum degree of $2n + 1$.

It is often useful to be able to fix one of the endpoints of the interval as one of the quadrature points. As an example, suppose we prescribe that $x_0 = a$. Let p be an arbitrary polynomial of degree $2n$, and write

$$p(x) = (x - a)q(x) + r,$$

where the quotient q is a polynomial of degree $2n - 1$ and the remainder r is a constant. The integral of $w p$ is then

$$\int_a^b w(x)p(x)dx = \int_a^b (x - a)w(x)q(x)dx + r \int_a^b w(x)dx.$$

We can now construct the usual Gauss quadrature formula for the interval $[a, b]$ with the modified weight function $(x - a)w(x)$, giving n quadrature points and n weights $x_k, W_k, k = 1, \dots, n$. This formula will be exact for all polynomials q of degree $2n - 1$. Provided that the weight function w satisfies the standard conditions on (a, b) , the modified weight function does also; in particular it is clearly positive on (a, b) . This gives

$$\begin{aligned} \int_a^b w(x)p(x)dx &= \sum_{k=1}^n W_k q(x_k) + r \int_a^b w(x)dx \\ &= \sum_{k=1}^n \frac{W_k}{x_k - a} p(x_k) \\ &\quad + r \int_a^b w(x)dx = \sum_{k=1}^n \frac{W_k}{x_k - a} p(x_k). \end{aligned} \tag{10.26}$$

The fact that $r = p$ (a) then leads us to consider the quadrature rule

$$\int_a^b w(x)f(x)dx \approx W f(a) + \sum_{k=1}^n W_k f(x_k), \quad (10.27)$$

where

$$\begin{aligned} W_k &= W / (x_k - a), \quad k = 1, \dots, n, \\ W &= \int_a^b w(x)dx. \end{aligned} \quad (10.28)$$

By construction, this formula is exact for all polynomials of degree $2n$. It is obvious that $W_k > 0$ for $k = 1, \dots, n$. We leave it as an exercise to show that $W > 0$ also (see Exercise 5).

With only trivial changes it is easy to see how to construct a similar formula where instead of fixing $x = a$ we fix $x = b$. These are known as **Radau quadrature formulae**. We leave it as an exercise to construct the formula corresponding to fixing both $x = a$ and $x = b$, which is known as a **Lobatto quadrature formula**; as might be expected, this is exact for all polynomials of degree $2n - 1$ (see Exercise 7).

The formal process could evidently be generalised to allow for fixing one of the quadrature points at an internal point c , where $a < c < b$. However, this leads to the difficulty that the modified weight function

$$w : x \rightarrow (x - c)w(x)$$

is not positive over the whole interval (a, b) ; hence we can no longer be sure that it is possible to construct a system of orthogonal polynomials, or, even if we can, that these polynomials will have all their zeros real and distinct and lying in $[a, b]$. In general, therefore, such quadrature formulae may not exist.

10.7 Note

For a detailed guide to the literature on Gauss quadrature rules and its connection to the theory of orthogonal polynomials, we refer to the books cited in the Notes at the end of Chapter 7.

Exercises

- 10.1 Determine the quadrature points and weights for the weight function $w: x \rightarrow \ln x$ on the interval $(0, 1)$, for $n = 0$ and $n = 1$.

- 10.2 The weights in the Gauss quadrature formula are given by (10.7), which is

$$W_k = \frac{1}{n!} w(x_k) [L'(x_k)]^{-1} dx.$$

Show that W_k can also be calculated from

$$W_k = \frac{1}{n!} w(x_k) L''(x_k).$$

(This is a simpler way of calculating W_k than (10.7); the importance of (10.7) is that it shows that the weights are all positive.)

- 10.3 Suppose that f has a continuous second derivative on $[0, 1]$. Show that there is a point ξ in $(0, 1)$ such that

$$\int_0^1 x f(x) dx = -f(\xi) + \frac{1}{2} f'(0).$$

- 10.4 Let $n \geq 0$. Write down the quadrature points $x_j, j = 0, \dots, n$, for the weight function $w(x) = (1-x)^{-1}$ on the interval $(-1, 1)$.

By induction, or otherwise, show that for positive integer values of n ,

$$\cos(2j+1) = \frac{\sin(2n+2)}{2 \sin 2},$$

unless n is a multiple of j . What is the value of the sum when n is a multiple of j ?

Deduce that

$$T_{n+1}(x) = \frac{1}{n} \int_{-1}^1 (1-x)^{-1} T_n(x) dx, \quad k = 1, \dots, n,$$

and show that

$$T_{n+1}(x) = \frac{n+1}{n} \int_{-1}^1 (1-x)^{-1} T_n(x) dx,$$

where T_n is the Chebyshev polynomial of degree n .

Deduce that the weights of the quadrature formula with weight function $w(x) = (1-x)^{-1}$ on the interval $(-1, 1)$ are

$$W_k = \frac{1}{n+1}, \quad k = 0, \dots, n.$$

- 10.5 In the notation for the construction of the Radau quadrature formula in Section 10.6, show that $W_n > 0$.

- 10.6 The **Laguerre polynomials** L_j , $j = 0, 1, 2, \dots$, are the orthogonal polynomials associated with the weight function $w(x) = e^{-x}$ on the semi-infinite interval $(0, \infty)$, with L_j of exact degree j . (See Exercise 5.9.) Show that

$$\int_0^\infty e^{-x} [L_j(x) - L_j(x)] p(x) dx = 0$$

when p is any polynomial of degree less than j .

In the Radau formula

$$\int_0^\infty e^{-x} p(x) dx = W_0 p(0) + \sum_{k=1}^n W_k p(x_k),$$

where one of the quadrature points is fixed at $x = 0$, show that the other quadrature points x_k , $k = 1, \dots, n$, are the zeros of the polynomial $L_{n+1} - L_n$. Deduce that

$$\int_0^\infty e^{-x} p(x) dx = -p'(0) + \sum_{k=1}^n W_k p(x_k).$$

- 10.7 Let $n \geq 2$. Show that a polynomial p_{2n-1} of degree $2n - 1$ can be written

$$p_{2n-1}(x) = (x - a)(b - x)q_{2n-3}(x) + r(x - a) + s(b - x),$$

where q_{2n-3} is a polynomial of degree $2n - 3$, and r and s are constants. Hence construct the Lobatto quadrature formula

$$\int_a^b w(x)f(x)dx = W_0 f(a) + \sum_{k=1}^{n-2} W_k f(x_k) + W_n f(b),$$

which is exact when f is any polynomial of degree $2n - 1$. Show that all the weights W_k , $k = 0, 1, \dots, n$, are positive.

- 10.8 Construct the Lobatto quadrature formula

$$\int_{-1}^1 f(x) dx = A_0 f(-1) + \sum_{k=1}^{n-2} A_k f(x_k) + A_n f(1)$$

for the interval $(-1, 1)$ with weight function $w(x) = 1$, and with $n = 2$; write down and solve four equations to determine x_1, A_1, A_0 and A_2 .

¹ \$/ , 9 , 40 : B@ % #.* \$ \$/ , 9 , C# % #..? \$

- 10.9 Write T for the composite trapezium rule (7.15), S for the composite Simpson rule (7.17) and M for the composite mid-point rule (10.25), each with m subintervals. Show that

$$M = 2I - T, \quad S = \frac{4I - T}{3}, \quad S = \frac{2M + I}{3}.$$

Piecewise polynomial approximation

11.1 Introduction

Up to now, the focus of our discussion has been the question of approximation of a given function f , defined on an interval $[a, b]$, by a polynomial on that interval either through Lagrange interpolation or Hermite interpolation, or by seeking the polynomial of best approximation (in the ∞ -norm or 2-norm). Each of these constructions was *global* in nature, in the sense that the approximation was defined by the same analytical expression on the whole interval $[a, b]$. An alternative and more flexible way of approximating a function f is to divide the interval $[a, b]$ into a number of subintervals and to look for a piecewise approximation by polynomials of low degree. Such piecewise-polynomial approximations are called **splines**, and the endpoints of the subintervals are known as the **knots**.

More specifically, a spline of degree n , $n \geq 1$, is a function which is a polynomial of degree n or less in each subinterval and has a prescribed degree of smoothness. We shall expect the spline to be at least continuous, and usually also to have continuous derivatives of order up to k for some k , $0 \leq k < n$. Clearly, if we require the derivative of order n to be continuous everywhere the spline is just a single polynomial, since if two polynomials have the same value and the same derivatives of every order up to n at a knot, then they must be the same polynomial. An important class of splines have degree n , with continuous derivatives of order up to and including $n - 1$, but as we shall see later, lower degrees of smoothness are sometimes considered.

To give a flavour of the theory of splines, we concentrate here on two simple cases: linear splines and cubic splines.

11.2 Linear interpolating splines

Definition 11.1 Suppose that f is a real-valued function, defined and continuous on the closed interval $[a, b]$. Further, let $K = x_0, \dots, x_m$ be a subset of $[a, b]$, with $a = x_0 < x_1 < \dots < x_m = b$, $m \geq 2$. The **linear spline** $S_{\#}$, interpolating f at the points x_i , is defined by

$$S_{\#}(x) = \frac{x - x_{i-1}}{x_i - x_{i-1}} f(x_{i-1}) + \frac{x_i - x}{x_i - x_{i-1}} f(x_i), \quad x \in [x_{i-1}, x_i], \quad i = 1, 2, \dots, m. \quad (11.1)$$

The points x_i , $i = 0, 1, \dots, m$, are the **knots** of the spline, and K is referred to as the **set of knots**.

As the function $S_{\#}$ interpolates the function f at the knots, i.e., $S_{\#}(x_i) = f(x_i)$, $i = 0, 1, \dots, m$, and over each interval $[x_{i-1}, x_i]$, for $i = 1, \dots, m$, the function $S_{\#}$ is a linear polynomial (and therefore continuous), we conclude that $S_{\#}$ is a continuous piecewise linear function on the interval $[a, b]$.

Given a set of knots $K = x_0, \dots, x_m$, we shall use the notation $h_i = x_i - x_{i-1}$, and let $h = \max h_i$. Also, for a positive integer k , we denote by $C^k[a, b]$ the set of all real-valued functions, defined and continuous on the closed interval $[a, b]$, such that all derivatives, up to and including order k , are defined and continuous on $[a, b]$.

In order to highlight the accuracy of interpolation by linear splines we state the following error bound in the ∞ -norm over the interval $[a, b]$.

Theorem 11.1 Suppose that $f \in C^2[a, b]$ and let $S_{\#}$ be the linear spline that interpolates f at the knots $a = x_0 < x_1 < \dots < x_m = b$; then, the following error bound holds:

$$\|f - S_{\#}\|_{\infty} \leq \frac{1}{8} h^2 \|f''\|_{\infty},$$

where $h = \max h_i = \max (x_i - x_{i-1})$, and $\| \cdot \|_{\infty}$ denotes the ∞ -norm over $[a, b]$, defined in (8.1).

Proof Consider a subinterval $[x_{i-1}, x_i]$, $1 \leq i \leq m$. According to Theorem 6.2, applied on the interval $[x_{i-1}, x_i]$,

$$f(x) - S_{\#}(x) = \frac{1}{2} f''(\xi)(x - x_{i-1})(x - x_i), \quad x \in [x_{i-1}, x_i],$$

where $h_i = (x_i - x_{i-1})$. Thus,

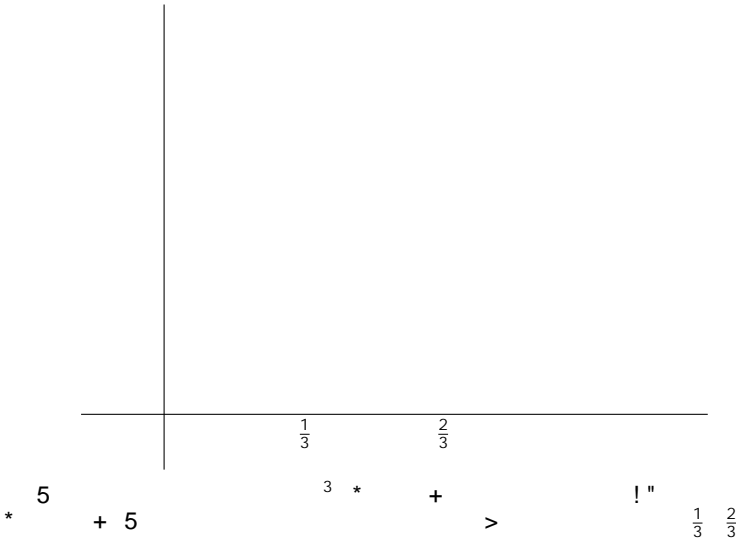
$$|f(x) - s_{\#}(x)| \leq \frac{1}{8} h_i \max_{x \in [x_{i-1}, x_i]} |f''(x)|.$$

Hence,

$$|f(x) - s_{\#}(x)| \leq \frac{1}{8} h \max_{x \in [a, b]} |f''(x)|,$$

for each $x \in [x_{i-1}, x_i]$ and each $i = 1, 2, \dots, m$. This gives the required error bound. □

Figure 11.1 shows a typical example: a linear spline approximation to the function $f: x \mapsto e^{-x}$ over the interval $[0, 1]$, using two internal knots, $x_1 = \frac{1}{3}$, $x_2 = \frac{2}{3}$, together with the endpoints of the interval, $x_0 = 0$ and $x_3 = 1$.



We conclude this section with a result that provides a characterisation of linear splines from the viewpoint of the calculus of variations.

A subset A of the real line is said to have **measure zero** if it can be contained in a countable union of open intervals of arbitrarily small total length; in other words, for every $\epsilon > 0$ there exists a sequence of open intervals (a_i, b_i) , $i = 1, 2, 3, \dots$, such that

$$A \subseteq \bigcup_{i=1}^{\infty} (a_i, b_i) \quad \text{and} \quad \sum_{i=1}^{\infty} (b_i - a_i) < \epsilon.$$

In particular, any finite or countable set A has measure zero. For example, the set of all rational numbers is countable, and therefore it has measure zero. Trivially, the empty set has measure zero.

Suppose that B is a subset of \mathbb{R} . We shall say that a certain property $P(x)$ holds for **almost every** x in B , if there exists a set $A \subset B$ of measure zero such that $P(x)$ holds for *all* $x \in B \setminus A$.

A real-valued function v defined on the interval $[a, b]$ is said to be **absolutely continuous** on $[a, b]$ if it has finite derivative $v'(x)$ at almost every point in $[a, b]$, v is (Lebesgue-) integrable on $[a, b]$, and

$$\int_a^x v'(t) dt = v(x) - v(a), \quad a \leq x \leq b.$$

Example 11.1 Any $v \in C^1[a, b]$ is absolutely continuous on the interval $[a, b]$. The function $x \mapsto x - (a+b)$ is absolutely continuous on $[a, b]$, but it does not belong to $C^1[a, b]$ as it is not differentiable at $x = -(a+b)$.

Let us denote by $H^1(a, b)$ the set of all absolutely continuous functions v defined on $[a, b]$ such that $v' \in L^1(a, b)$, i.e.,

$$\int_a^b |v'(x)| dx < \infty.$$

We observe in passing that any function $v \in H^1(a, b)$ is uniformly continuous on the closed interval $[a, b]$. This follows by noting that, for any pair of points $x, y \in [a, b]$,

$$\begin{aligned} |v(x) - v(y)| &= \left| \int_x^y v'(t) dt \right| \\ &\leq \int_x^y |v'(t)| dt \\ &\leq \int_x^y |v'(t)|^k dt^{1/k} \cdot (y-x)^{(k-1)/k} \\ &\leq \int_x^y |v'(t)|^k dt^{1/k} \cdot |y-x|^{(k-1)/k}. \end{aligned}$$

In the transition from the first line to the second we used the Cauchy-Schwarz inequality.

If $k \geq 1$, we shall denote by $H^k(a, b)$ the set of all $v \in H^1(a, b)$ such that $v^{(k)}$ is absolutely continuous on $[a, b]$ and $v^{(k)} \in L^1(a, b)$. The set $H^k(a, b)$ is called a **Sobolev space** of index k . We observe that

$$C^k[a, b] \subset H^k(a, b)$$

for any $k \geq 1$, with strict inclusion. For example, any linear spline on

$[a, b]$ belongs to $H(a, b)$, but not to $C[a, b]$ unless it is a linear function over the *whole* of the interval $[a, b]$.

Example 11.2 Let $\alpha > 1/2$; the function x^α then belongs to $H(0, 1)$, although it only belongs to $C[0, 1]$ if $\alpha \geq 1$.

As a second example, consider the function $x \ln x$ which belongs to $H(0, 1)$, but not to $C[0, 1]$.

The variational characterisation of linear splines stated in the next theorem expresses the fact that, among all functions $v \in H(a, b)$ which interpolate a given continuous function f at a fixed set of knots in $[a, b]$, the linear spline $S_\#$ that interpolates f at these knots is the ‘flattest’, in the sense that its ‘average slope’ $S_\#$ is smallest.

Theorem 11.2 Suppose that $S_\#$ is the linear spline that interpolates $f \in C[a, b]$ at the knots $a = x_0 < x_1 < \dots < x_m = b$. Then, for any function v in $H(a, b)$ that also interpolates f at these knots,

$$S_\# \leq v.$$

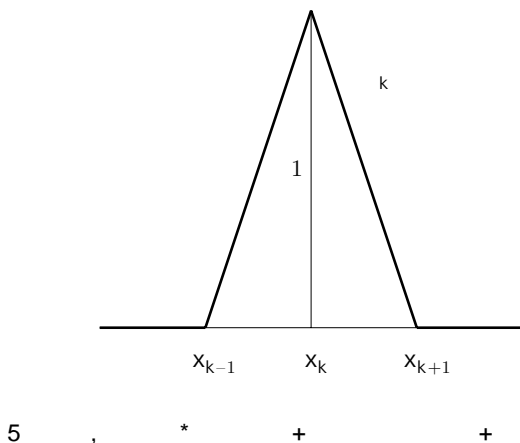
Proof Let us observe that

$$\begin{aligned} \int_a^b (v(x) - S_\#(x))^2 dx &= \int_a^b (v(x) - S_\#(x)) dx + \int_a^b S_\#(x) dx \\ &\quad + 2 \int_a^b (v(x) - S_\#(x))S_\#(x) dx. \end{aligned} \tag{11.2}$$

We shall now use integration by parts to show that the last integral is equal to 0; the desired inequality will then follow by noting that the first term on the right-hand side is nonnegative and it is equal to 0 if, and only if, $v = S_\#$. Clearly,

$$\begin{aligned} \int_a^b (v(x) - S_\#(x))S_\#(x) dx &= \int_a^b (v(x) - S_\#(x))S_\#(x) dx \\ &= \left[(v(x) - S_\#(x))S_\#(x) \right]_a^b - \int_a^b (v'(x) - S_\#'(x))S_\#(x) dx \\ &\quad - \int_a^b (v(x) - S_\#(x))S_\#'(x) dx. \end{aligned} \tag{11.3}$$

Now $v(x) - S_\#(x) = f(x) - f(x) = 0$ for $i = 0, 1, \dots, m$ and, since $S_\#$ is a linear polynomial over each of the open intervals (x_{k-1}, x_k) , $k =$



$1, 2, \dots, m$, it follows that $S_{\#}$ is identically 0 on each of these intervals. Thus, the expression in the square bracket in (11.3) is equal to 0 for each $k = 1, 2, \dots, m$. \square

Sobolev spaces play an important role in approximation theory. We shall encounter them again in Chapter 14 which is devoted to the approximation of solutions to differential equations by piecewise polynomial functions.

11.3 Basis functions for the linear spline

Suppose that $S_{\#}$ is a linear spline with knots $x_i, i = 0, 1, \dots, m$, interpolating the function $f \in C[a, b]$. Instead of specifying the value of $S_{\#}$ on each subinterval $[x_{i-1}, x_i], i = 1, 2, \dots, m$, we can express $S_{\#}$ as a linear combination of suitable ‘basis functions’ as follows:

$$S_{\#}(x) = \sum_{i=1}^m \phi_i(x) f(x_i), \quad x \in [a, b].$$

Here, we require that each ϕ_i is itself a linear spline which vanishes at every knot except x_i , and $\phi_i(x_i) = 1$. The function ϕ_i is often known as the **linear basis spline** or **hat function**, and is depicted in Figure 11.2.

The formal definition of $s_k(x)$ is as follows:

$$s_k(x) = \begin{cases} 0 & \text{if } x = x_{k-1}, \\ (x - x_{k-1})/h & \text{if } x_{k-1} < x < x_k, \\ (x_k - x)/h & \text{if } x_k < x < x_{k+1}, \\ 0 & \text{if } x = x_{k+1}, \end{cases}$$

for $k = 1, \dots, m - 1$, and with

$$s_0(x) = \begin{cases} (x - a)/h & \text{if } a = x < x_0 < x_1, \\ 0 & \text{if } x = x_0 \end{cases} \tag{11.4}$$

and

$$s_m(x) = \begin{cases} 0 & \text{if } x = x_{m-1}, \\ (x - x_{m-1})/h & \text{if } x_{m-1} < x < x_m = b. \end{cases}$$

11.4 Cubic splines

Suppose that $f \in C[a, b]$ and let $K = \{x_0, \dots, x_m\}$ be a set of $m + 1$ knots in the interval $[a, b]$, $a = x_0 < x_1 < \dots < x_m = b$. Consider the set \mathcal{S} of all functions $s \in C[a, b]$ such that

$$\begin{aligned} s(x) &= f(x), \quad i = 0, 1, \dots, m, \\ s &\text{ is a cubic polynomial on } [x_{i-1}, x_i], \quad i = 1, 2, \dots, m. \end{aligned}$$

Any element of \mathcal{S} is referred to as an **interpolating cubic spline**. We note that, unlike linear splines which are uniquely determined by the interpolating conditions, there is more than one interpolating cubic spline $s \in C[a, b]$ that satisfies the two conditions stated above; indeed, there are $4m$ coefficients of cubic polynomials (four on each subinterval $[x_{i-1}, x_i]$, $i = 1, 2, \dots, m$), and only $m + 1$ interpolating conditions and $3(m - 1)$ continuity conditions; since s belongs to $C[a, b]$, this means that s , s' and s'' are continuous at the internal knots x_1, \dots, x_{m-1} . Hence, we have a total of $4m - 2$ conditions for the $4m$ unknown coefficients. Depending on the choice of the remaining two conditions we can construct various interpolating cubic splines.

An important class of cubic splines is singled out by the following definition.

Definition 11.2 *The natural cubic spline, denoted by s_n , is the element of the set \mathcal{S} satisfying the end conditions*

$$s_n'(x_0) = s_n'(x_m) = 0.$$

We shall prove that this definition is correct in the sense that the two additional conditions in Definition 11.2 uniquely determine S : this will be done by describing an algorithm for constructing S .

Construction of the natural cubic spline. Let us begin by defining $s_i = S(x_i)$, $i = 0, 1, \dots, m$, and noting that S is a linear function on each subinterval $[x_{i-1}, x_i]$. Therefore, S can be expressed as

$$s_i(x) = \frac{x - x_{i-1}}{h_i} s_{i-1} + \frac{x_i - x}{h_i} s_i$$

Theorem 11.3 Let s be the natural cubic spline that interpolates a function $f \in C[a, b]$ at the knots $a = x_0 < x_1 < \dots < x_m = b$. Then, for any function v in $H^1(a, b)$ that also interpolates f at the knots,

$$\int_a^b |s'(x) - v'(x)|^2 dx \leq \int_a^b |v'(x)|^2 dx.$$

The proof is analogous to that of Theorem 11.2 and is left as an exercise.

The *smoothest interpolation property* expressed by Theorem 11.3 is the source of the name *spline*. A spline is a flexible thin curve-drawing aid, made of wood, metal or acrylic. Assuming that its shape is given by the equation $y = v(x)$, $x \in [a, b]$, and is constrained by requiring that it passes through a finite set of prescribed points in the plane, v will take on a shape which minimises the strain energy

$$E(v) = \int_a^b \frac{1}{2} |v'(x)|^2 dx$$

over all functions v which are constrained in the same way. If the function v is slowly varying, *i.e.*, $\max |v'(x)| \ll 1$, this energy-minimisation property is very similar to the result in Theorem 11.3.

11.5 Hermite cubic splines

In the previous section we took $f \in C[a, b]$ and demanded that s belonged to $C^1[a, b]$; here we shall strengthen our requirements on the smoothness of the function that we wish to interpolate and assume that $f \in C^1[a, b]$; simultaneously, we shall relax the smoothness requirements on the associated spline approximation s by demanding that $s \in C^0[a, b]$ only.

Let $K = \{x_0, \dots, x_m\}$ be a set of knots in the interval $[a, b]$ with $a = x_0 < x_1 < \dots < x_m = b$ and $m \geq 2$. We define the **Hermite cubic spline** as a function $s \in C^0[a, b]$ such that

$$\begin{aligned} s(x) &= f(x), \quad s'(x) = f'(x) \text{ for } i = 0, 1, \dots, m, \\ s &\text{ is a cubic polynomial on } [x_{i-1}, x_i] \text{ for } i = 1, 2, \dots, m. \end{aligned}$$

Writing the spline s on the interval $[x_{i-1}, x_i]$ as

$$s(x) = c_0 + c_1(x - x_{i-1}) + c_2(x - x_{i-1})^2 + c_3(x - x_{i-1})^3, \quad x \in [x_{i-1}, x_i], \quad (11.8)$$

¹ $\int_a^b |v'(x)|^2 dx$ is the strain energy of a beam of length $b - a$ fixed at $x = a$ and free at $x = b$. ² $\int_a^b |v'(x)|^2 dx$ is the strain energy of a beam of length $b - a$ fixed at both ends. ³ $\int_a^b |v'(x)|^2 dx$ is the strain energy of a beam of length $b - a$ fixed at both ends and with a point load at $x = b/2$. ⁴ $\int_a^b |v'(x)|^2 dx$ is the strain energy of a beam of length $b - a$ fixed at both ends and with a point load at $x = b/2$.

we find that $c_0 = f(x_{i-1})$, $c_1 = f'(x_{i-1})$, and

$$\begin{aligned} c_2 &= 3 \frac{f(x_i) - f(x_{i-1})}{h} - \frac{f'(x_i) + 2f'(x_{i-1}))}{h}, \\ c_3 &= \frac{f(x_i) + f(x_{i-1}))}{2h} - \frac{f'(x_i) - f'(x_{i-1}))}{h}. \end{aligned} \tag{11.9}$$

Note that the Hermite cubic spline only has a continuous first derivative at the knots, and therefore it is *not* an interpolating cubic spline in the sense of Section 11.4.

Unlike natural cubic splines, the coefficients of a Hermite cubic spline on each subinterval can be written down explicitly without the need to solve a tridiagonal system.

Concerning the size of the interpolation error, we have the following result.

Theorem 11.4 *Let $f \in C^4[a, b]$, and let s be the Hermite cubic spline that interpolates f at the knots $a = x_0 < x_1 < \dots < x_n = b$; then, the following error bound holds:*

$$\|f - s\| \leq \frac{1}{384} h^4 \|f^{(4)}\|,$$

where $f^{(4)} = f^{(4)}$ is the fourth derivative of f with respect to its argument, x , $h = \max_{i=1, \dots, n} h_i = \max_{i=1, \dots, n} (x_i - x_{i-1})$, and $\|\cdot\|$ denotes the ∞ -norm on the interval $[a, b]$.

The proof is analogous to that of Theorem 11.1, except that Theorem 6.4 is used instead of Theorem 6.2.

Both the linear spline and the Hermite cubic spline are local approximations; the value of the spline at a point x between two knots x_{i-1} and x_i depends only on the values of the function and its derivative at these two knots. On the other hand, the natural cubic interpolating spline is a global approximation and, in this respect, it is more typical of a generic spline: a change in just one of the values at a knot, $f(x_i)$, will alter the right-hand side of the system of equations (11.7), so the values of all the quantities c_j will change. Thus, the spline will change throughout the whole interval $[x_{i-1}, x_i]$. We conclude this section with an example.

Example 11.3 *Figure 11.3 shows the Hermite cubic spline approximation to the function $f: x \mapsto 1/(1+x)$, using four equally spaced knots in the interval $[0, 5]$.*



from the standard properties of binomial coefficients. Hence Q_{r-1} is a polynomial in x of degree $n - r + 1$, and the result follows by induction. Finally, this shows that Q_n is a polynomial of degree 0, and is therefore constant on I_n . Thus, by the same argument, Q_{n-1} is identically 0 on I_n . \square

Theorem 11.5 For each $n \geq 1$, the function S_n defined by

$$S_n(x) = \binom{n}{k} \frac{(x - kh)^{n-k}}{k!} (x - kh)$$

is a spline of degree n with equally spaced knots kh , $k = 0, 1, \dots, n + 1$. It has a continuous derivative of order $n - 1$ and is identically 0 outside the interval $(0, (n + 1)h)$.

Proof The function S_n is clearly a spline as stated, and $S_n(x)$ is identically 0 for $x \leq 0$. When $x \in (n + 1)h$ the arguments $x - kh$, $k = 0, 1, \dots, n + 1$, of the positive parts are all nonnegative, so that

$$S_n(x) = \binom{n}{k} \frac{(x - kh)^{n-k}}{k!} (x - kh),$$

and this is identically zero by Lemma 11.1. \square

Taking $n = 1$ we find that

$$S_1(x) = x - 2(x - h) + (x - 2h).$$

After normalisation by $1/h$ so as to have a maximum value of 1, and shifting $x = 0$ to $x = x -$, this yields a representation of the linear hat function ϕ_1 from (11.4) in the form

$$\phi_1(x) = \frac{1}{h} S_1(x - x_-),$$

which, for $1 - k \leq x \leq k$, is nonzero over two consecutive intervals: $(x - , x -]$ and $[x - , x -)$.

In the same way we obtain a basis function for the cubic spline by taking $n = 3$:

$$S_3(x) = x - 4(x - h) + 6(x - 2h) - 4(x - 3h) + (x - 4h).$$

Normalising so as to have a maximum value of 1 and shifting $x = 0$ to $x = x -$, we get

$$\phi_3(x) = \frac{1}{4h} S_3(x - x_-).$$

proximation in Theorem 11.4, but not for the natural cubic spline. The analysis of the error in the natural cubic spline approximation is quite complicated; Powell gives full details in his book.

The following are classical texts on the theory of splines.

De Boor, C. A., *A Practical Guide to Splines*, Mathematics in Science and Engineering, 38, Academic Press, New York, 1967.

De Boor, C. A., *A Practical Guide to Splines*, Revised Edition, Springer Applied Mathematical Sciences, 27, Springer, New York, 2001.

De Boor, C. A. & Ronken, K., *Spline Functions: Basic Theory*, John Wiley & Sons, New York, 1981.

The variational characterisations of splines stated in Sections 11.1 and 11.3 stem from the work of J.C. Holladay, Smoothest curve approximation, *Math. Comput.* **11**, 233–243, 1957.

Our definition of the Sobolev space $H^1(a, b)$ in Section 11.1, based on the concept of absolute continuity, is specific to functions of a single variable. More generally, for functions of several real variables one needs to invoke the theory of weak differentiability or the theory of distributions to give a rigorous definition of the Sobolev space $H^1(\Omega)$ with $\Omega \subset \mathbb{R}^n$; alternatively, one can define $H^1(\Omega)$ by completion of the set of smooth functions in a suitable norm. For the sake of simplicity of exposition we have chosen to avoid such general approaches.

Exercises

- 11.1 An interpolating spline of degree n is required to have continuous derivatives of order up to and including $n - 1$ at the knots. How many additional conditions are required to specify the spline uniquely?
- 11.2 (i) Suppose that f is a polynomial of degree 1. Show that the linear spline $S_{\#}$ which interpolates f at the knots x_i for $i = 0, 1, \dots, m$ is identical to f , so that $S_{\#} = f$.
 (ii) Suppose that f is a polynomial of degree 3. Show that the Hermite cubic spline S_{\S} which interpolates f at the knots x_i , $i = 0, 1, \dots, m$, is identical to f , so that $S_{\S} = f$.
 (iii) Suppose that f is a polynomial of degree 3. Show that the natural cubic spline S which interpolates f at the knots x_i , $i = 0, 1, \dots, m$, is not in general identical to f .

- 11.3 Suppose that the natural cubic spline S interpolates the function $f: x \rightarrow x$ on the interval $[0, 1]$, the knots being equally spaced, so that $x_i = ih$, $i = 0, 1, \dots, m$, with $h = 1/m$, $m \geq 2$. Write down the equations which determine the quantities c_i . If the two additional conditions are $c_0 = c_m = 0$, show that these equations are not satisfied by $c_i = f'(x_i)$, $i = 1, \dots, m-1$, so that S and f are not identical. If, however, these two additional conditions are replaced by $c_0 = f'(0)$, $c_m = f'(1)$, show that $c_i = f'(x_i)$, $i = 0, 1, \dots, m$, and deduce that S and f are identical.
- 11.4 A linear spline on the interval $[0, 1]$ is expressed in terms of the basis functions as

$$S(x) = \sum_{i=0}^{m-1} c_i \phi_i(x).$$

Instead of being required to interpolate the function f at the knots, the spline S is required to minimise $\int_0^1 (f - S)^2 dx$. Show that the coefficients c_i satisfy the system of equations

$$A c = b,$$

where the elements of the matrix A are

$$A_{ij} = \int_0^1 \phi_i(x) \phi_j(x) dx$$

and the elements of b are

$$b_i = \int_0^1 f(x) \phi_i(x) dx.$$

Now suppose that the knots are equally spaced, so that $x_k = kh$, $k = 0, 1, \dots, m$, where $h = 1/m$, $m \geq 2$. Show that the matrix A is tridiagonal, with $A_{ii} = -h$ for $i = 1, \dots, m-1$, and determine the other nonzero elements of A . Show also that A has the properties required for the use of the Thomas algorithm described in Section 3.3.

- 11.5 In the notation of Exercise 4, suppose that $f(x) = x$. Verify that the system of equations is satisfied by $c_i = kh$, so that $S = f$.

Now suppose that $f(x) = x^2$. Verify that the equations are satisfied by $c_i = (kh)^2 + Ch$, where C is a constant to be determined. Deduce that $S(x) = f(x) + Ch$.

- 11.6 In the notation of Theorem 11.5, the spline basis function S of degree n is defined by

$$S(x) = \binom{n+1}{k} (x - kh)^k.$$

Explain why, for any value of a ,

$$(x - a)^k (x - a) = (x - a)^{k+1}.$$

Show that

$$xS(x) + [(n+2)h - x]S(x - h) = S(x).$$

Hence show by induction that $S(x) \geq 0$ for all x .

- 11.7 Use the result of Exercise 6 to show by induction that each basis function S is symmetric; that is,

$$S(p + x) = S(p - x)$$

for all x , where $p = -(n+1)h$.

Initial value problems for ODEs

12.1 Introduction

Ordinary differential equations frequently occur in mathematical models that arise in many branches of science, engineering and economics. Unfortunately it is seldom that these equations have solutions which can be expressed in closed form, so it is common to seek approximate solutions by means of numerical methods. Nowadays this can usually be achieved very inexpensively to high accuracy and with a reliable bound on the error between the analytical solution and its numerical approximation. In this section we shall be concerned with the construction and the analysis of numerical methods for first-order differential equations of the form

$$y' = f(x, y) \tag{12.1}$$

for the real-valued function y of the real variable x , where $y' = \frac{dy}{dx}$ and f is a given real-valued function of two real variables. In order to select a particular integral from the infinite family of solution curves that constitute the general solution to (12.1), the differential equation will be considered in tandem with an **initial condition**: given two real numbers x_0 and y_0 , we seek a solution to (12.1) for $x > x_0$ such that

$$y(x_0) = y_0. \tag{12.2}$$

The differential equation (12.1) together with the initial condition (12.2) is called an **initial value problem**.

If you believe that any initial value problem of the form (12.1), (12.2) possesses a unique solution, take a look at the following example.

Proof We define a sequence of functions (y_n) by

$$y_0(x) = y_0, \quad y_n(x) = y_0 + \int_{x_0}^x f(s, y_{n-1}(s)) ds, \quad n = 1, 2, \dots \quad (12.4)$$

Since f is continuous on D , it is clear that each function y_n is continuous on $[x_0, X_0]$. Further, since

$$y_n(x) = y_0 + \int_{x_0}^x f(s, y_n(s)) ds,$$

it follows by subtraction that

$$y_n(x) - y_{n-1}(x) = \int_{x_0}^x [f(s, y_n(s)) - f(s, y_{n-1}(s))] ds. \quad (12.5)$$

We now proceed by induction, and assume that, for some positive value of n ,

$$|y_n(x) - y_{n-1}(x)| \leq \frac{K [L(x - x_0)]^n}{n!}, \quad x \in [x_0, X_0], \quad (12.6)$$

and that

$$|y_k(x) - y_{k-1}(x)| \leq \frac{K [L(x - x_0)]^k}{k!}, \quad x \in [x_0, X_0], \quad k = 1, \dots, n. \quad (12.7)$$

Trivially, the hypotheses of the theorem and (12.4) imply that (12.6) and (12.7) hold for $n = 1$.

Now, (12.7) and (12.3) yield that

$$|y_k(x) - y_{k-1}(x)| \leq \frac{K}{L} e^{L(x - x_0)} - 1 = C, \quad x \in [x_0, X_0], \quad k = 1, \dots, n.$$

Therefore $(x, y_n(x)) \in D$ and $(x, y_{n-1}(x)) \in D$ for all $x \in [x_0, X_0]$. Hence, using (12.5), the Lipschitz condition and (12.6),

$$\begin{aligned} |y_n(x) - y_{n-1}(x)| &\leq L \int_{x_0}^x \frac{K [L(s - x_0)]^{n-1}}{(n-1)!} ds \\ &= \frac{K [L(x - x_0)]^n}{n!}, \end{aligned} \quad (12.8)$$

for all $x \in [x_0, X_0]$. Moreover, using (12.8) and (12.7),

$$\begin{aligned} y^{(j)}(x) - y^{(j)}(x_0) &= \frac{K}{L} \frac{[L(x - x_0)]^j}{(j+1)!} + \frac{K}{L} \frac{[L(x - x_0)]^j}{j!} \\ &= \frac{K}{L} \frac{[L(x - x_0)]^j}{(j+1)!}, \end{aligned} \quad (12.9)$$

for all $x \in [x_0, X_0]$. Thus, (12.6) and (12.7) hold with n replaced by $n+1$, and hence, by induction, they hold for all positive integers n .

Since the infinite series $\sum (c^j/j!)$ converges (to $e^c - 1$) for any value of c , and for $c = L(X_0 - x_0)$ in particular, it follows from (12.6) that the infinite series

$$[y^{(j)}(x) - y^{(j)}(x_0)]$$

converges absolutely and uniformly for $x \in [x_0, X_0]$. However,

$$y^{(j)}(x) + [y^{(j)}(x) - y^{(j)}(x_0)] = y^{(j)}(x_0),$$

showing that the sequence of continuous functions $(y^{(j)})$ converges to a limit, uniformly on $[x_0, X_0]$, and hence that the limit itself is a continuous function. Calling this limit y , we see from (12.4) that

$$\begin{aligned} y(x) &= \lim_{j \rightarrow \infty} y^{(j)}(x) \\ &= y(x_0) + \lim_{j \rightarrow \infty} \int_{x_0}^x f(s, y^{(j)}(s)) ds, \\ &= y(x_0) + \int_{x_0}^x \lim_{j \rightarrow \infty} f(s, y^{(j)}(s)) ds, \\ &= y(x_0) + \int_{x_0}^x f(s, y(s)) ds, \end{aligned} \quad (12.10)$$

where we used the uniform convergence of the sequence of functions $(y^{(j)})$ in the transition from line two to line three to interchange the order of the limit process and integration, and the continuity of the function f in the transition from line three to line four. As $s \mapsto f(s, y(s))$ is a continuous function of s on the interval $[x_0, X_0]$, its integral over the interval $[x_0, x]$ is a continuously differentiable function of x . Hence, by

(12.10), y is a continuously differentiable function of x on $[x_0, X_0]$; *i.e.*,
 y

As a very simple example, consider the linear equation

$$y' = py + q, \quad (12.12)$$

where p and q are constants. Then, $L = p$, independently of C , and $K = py + q$. Hence, for any interval $[x, X_0]$, the conditions are satisfied by choosing C sufficiently large; therefore, the initial value problem has a unique continuously differentiable solution, defined for all $x \in [x, X_0]$.

Now, consider another example

$$y' = y^2, \quad y(0) = 1.$$

Here for any interval $[0, X_0]$ we have $K = 1$. Choosing any positive value of C we find that

$$u' - v' = u + v \quad u' - v' = L u - v \quad u, v \in C,$$

where $L = 2(1 + C)$. We therefore now require the condition

$$C \geq \frac{1}{2(1+C)} e^{(X_0)^2} - 1.$$

This is satisfied if

$$X_0 \leq F(C) = \frac{1}{2(1+C)} \ln(1 + 2C + 2C^2),$$

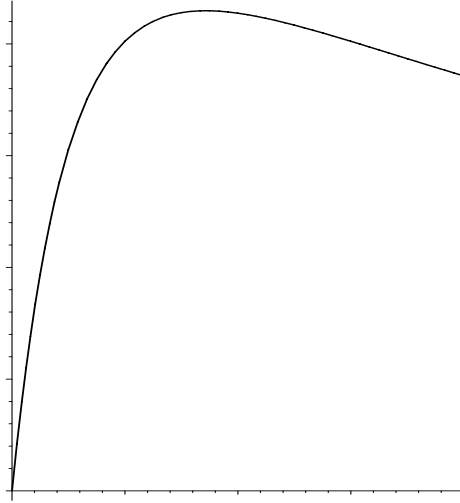
where \ln means \log_e . A sketch of the graph of the function F against C shows that F takes its maximum value near $C = 1.714$, and this gives the condition $X_0 \leq 0.43$ (see Figure 12.1).

Thus, we are *unable* to prove the existence of the solution over the infinite interval $[0, \infty)$. This is correct, of course, as the unique solution of the initial value problem is

$$y(x) = \frac{1}{1-x}, \quad 0 \leq x < 1,$$

and this is not continuous, let alone continuously differentiable, on any interval $[0, X_0]$ with $X_0 = 1$. The conditions of Picard's Theorem, which are sufficient but not necessary for the existence and the uniqueness of the solution, have given a rather more restrictive bound on the size of the interval over which the solution exists.

The method of proof of Picard's Theorem also suggests a possible technique for constructing approximations to the solution, by determining the functions y_n from (12.4). In practice it may be impossible, or very difficult, to evaluate the necessary integrals in closed form. We



! - #)' . *- + '& - ' .

leave it as an exercise (see Exercise 3) to show that for the simple linear equation (12.12), with initial condition $y(0) = 1$, the function y is the same as the approximation obtained from the exact solution by expanding the exponential function as a power series and retaining the terms up to the one involving x .

In the rest of this chapter we shall consider step-by-step numerical methods for the approximate solution of the initial value problem (12.1), (12.2). We shall suppose throughout that the function f satisfies the conditions of Picard's Theorem. Suppose that the initial value problem (12.1), (12.2) is to be solved on the interval $[x_0, X_0]$. We divide this interval by the **mesh points** $x_n = x_0 + nh$, $n = 0, 1, \dots, N$, where $h = (X_0 - x_0)/N$ and N is a positive integer. The positive real number h is called the **step size** or **mesh size**. For each n we seek a numerical approximation y_n to $y(x_n)$, the value of the analytical solution at the mesh point x_n ; these values y_n are calculated in succession, for $n = 1, 2, \dots, N$.

12.2 One-step methods

A one-step method expresses y_{n+1} in terms of the previous value y_n ; later on we shall consider k -step methods, where y_{n+1} is expressed in terms of the k previous values y_{n-k+1}, \dots, y_n , where $k \geq 2$. The simplest example of a one-step method for the numerical solution of the initial value problem (12.1), (12.2) is Euler's method.

Euler's method. Given that $y(x_0) = y_0$, let us suppose that we have already calculated y_n , up to some n , $0 \leq n \leq N-1$, $N \geq 1$; we define

$$y_{n+1} = y_n + hf(x_n, y_n).$$

Thus, taking in succession $n = 0, 1, \dots, N-1$, one step at a time, the approximate values y_n at the mesh points x_n can be easily obtained. This numerical method is known as **Euler's method**.

In order to motivate the definition of Euler's method, let us observe that on expanding $y(x_0 + h) = y(x_0 + h)$ into a Taylor series about x_0 , retaining only the first two terms, and writing $y(x_0 + h) = f(x_0, y(x_0))$, we have that

$$y(x_0 + h) = y(x_0) + hf(x_0, y(x_0)) + O(h^2).$$

After replacing $y(x_0)$ and $y(x_0 + h)$ by their numerical approximations, denoted by y_n and y_{n+1} , respectively, and discarding the $O(h^2)$ term, we arrive at Euler's method.

More generally, a one-step method may be written in the form

$$y_{n+1} = y_n + h\Phi(x_n, y_n; h), \quad n = 0, 1, \dots, N-1, \quad y(x_0) = y_0, \quad (12.13)$$

where $\Phi(x, y; h)$ is a continuous function of its variables. For example, in the case of Euler's method, $\Phi(x, y; h) = f(x, y)$. More intricate examples of one-step methods will be discussed below.

In order to assess the accuracy of the numerical method (12.13), we define the **global error**, e_n , by

$$e_n = y(x_n) - y_n.$$

We also need the concept of **truncation error**, τ_n , defined by

$$\tau_n = \frac{y(x_{n+1}) - y(x_n)}{h} - \Phi(x_n, y(x_n); h). \quad (12.14)$$

The next theorem provides a bound on the magnitude of the global error in terms of the truncation error.

Theorem 12.2

Assuming that $y \in C^2[x_0, X_0]$, i.e., that y is a twice continuously differentiable function of x on $[x_0, X_0]$, and expanding $y(x_0 + \tau)$ about the point x_0 into a Taylor series with remainder (see Theorem A.4), we have that

$$y(x_0 + \tau) = y(x_0) + \tau y'(x_0) + \frac{\tau^2}{2!} y''(\xi), \quad x_0 < \xi < x_0 + \tau.$$

Substituting this expansion into (12.18) gives

$$\tau = \frac{1}{2} \tau^2 y''(\xi).$$

Let $M = \max_{x \in [x_0, X_0]} |y''(x)|$. Then, $\tau = \frac{1}{2} \tau^2 M$, $n = 0, 1, \dots, N-1$, where $\tau = hM$. Inserting this into (12.16) and noting that for Euler's method $\Phi(x, y; h) = f(x, y)$ and therefore $L = 1$ where L is the Lipschitz constant for f , we have that

$$e_n \leq \frac{1}{2} M \frac{e^{1-L} - 1}{L} h, \quad n = 0, 1, \dots, N. \quad (12.19)$$

Let us highlight the practical relevance of our error analysis by focusing on a particular example.

Example 12.2 *Let us consider the initial value problem $y' = \tan^{-1} y$, $y(0) = y_0$, where y_0 is a given real number. In order to find an upper bound on the global error $e_N = y(x_N) - y_N$, where y_N is the Euler approximation to $y(x_N)$, we need to determine the constants L and M in the inequality (12.19).*

Here $f(x, y) = \tan^{-1} y$; so, by the Mean Value Theorem (Theorem A.3),

$$f(x, u) - f(x, v) = \frac{f'(x, \eta)}{1} (u - v) = \frac{1}{1 + \eta^2} (u - v),$$

where η lies between u and v . In our case

$$\frac{f'(x, y)}{1} = (1 + y^2)^{-1},$$

and therefore $L = 1$. To find M we need to obtain a bound on y' (without actually solving the initial value problem!). This is easily achieved by differentiating both sides of the differential equation with respect to the variable x :

$$y' = \frac{d}{dx}(\tan^{-1} y) = (1 + y^2)^{-1} \frac{dy}{dx} = (1 + y^2)^{-1} \tan^{-1} y.$$

Therefore $y(x) - M = -$. Inserting the values of L and M into (12.19) and noting that $x = 0$, we have

$$e - (e - 1)h, \quad n = 0, 1, \dots, N.$$

Thus, given a tolerance ϵ , specified beforehand, we can ensure that the error between the (unknown) analytical solution and its numerical approximation does not exceed this tolerance by choosing a positive step size h such that

$$h \leq \frac{\epsilon}{(e^2 - 1)}.$$

For such h we shall have $y(x) - y = \epsilon$, for $n = 0, 1, \dots, N$, as required. Thus, at least in principle, we can calculate the numerical solution to arbitrarily high accuracy by choosing a sufficiently small step size h .

A numerical experiment shows that this error estimate is rather pessimistic. Taking, for example, $y = 1$ and $X_0 = 1$, our bound implies that the tolerance $\epsilon = 0.01$ will be achieved with $h \approx 0.0074$; hence, it would appear that we need $N \approx 135$. In fact, using $N = 27$ gives a result from Euler's method which is just within this tolerance, so the error estimate has predicted the use of a step size which is five times smaller than is actually required.

Example 12.3 *As a more typical practical example, consider the problem*

$$y' = y + g(x), \quad y(0) = 2, \quad (12.20)$$

where

$$g(x) = \frac{x^2 - 6x + 12x - 14x + 9}{(1+x)},$$

is so chosen that the solution is known, and is

$$y(x) = \frac{(1-x)(2-x)}{1+x}.$$

The results of some numerical calculations on the interval $x \in [0, 1.6]$ are shown in Figure 12.2. They use step sizes 0.2, 0.1 and 0.05, and show how halving the step size gives a reduction of the error also by a factor of roughly 2, in agreement with the error bound (12.19).



that the one-step method (12.13) is consistent if, and only if,

$$\Phi(x, y; 0) = f(x, y). \tag{12.21}$$

This condition is sometimes taken as the definition of consistency. We shall henceforth always assume that (12.21) holds.

Now, we are ready to state a convergence theorem for the general one-step method (12.13).

Theorem 12.3 *Suppose that the initial value problem (12.1), (12.2) satisfies the conditions of Picard’s Theorem, and also that its approximation generated from (12.13) when $h = h_n$ lies in the region D . Assume further that the function $\Phi(x, y; h)$ is continuous on $D = [0, h_n]$, and satisfies the consistency condition (12.21) and the Lipschitz condition*

$$|\Phi(x, u; h) - \Phi(x, v; h)| \leq L |u - v| \quad \text{on } D = [0, h_n]. \tag{12.22}$$

Then, if successive approximation sequences $\{y_n(x)\}$, generated by using the mesh points $x_n = x_0 + nh_n$, $n = 1, 2, \dots, N$, are obtained from (12.13) with successively smaller values of h , each h_n less than h , we have convergence of the numerical solution to the solution of the initial value problem in the sense that

$$\lim_{n \rightarrow \infty} y_n(x) = y(x) \quad \text{as } x \rightarrow x_0 \text{ in } [x_0, X_0] \text{ when } h_n \rightarrow 0 \text{ and } n \rightarrow \infty.$$

Proof Suppose that $h = (X_0 - x_0)/N$, where N is a positive integer. We shall assume that N is sufficiently large so that $h < h$. Since $y(x_0) = y_0$ and therefore $e = 0$, Theorem 12.2 implies that

$$|y(x_n) - y_n(x_n)| \leq \frac{e^{L(x_n - x_0)} - 1}{L} \max_{x \in [x_0, x_n]} |\tau_n|, \quad n = 1, 2, \dots, N. \tag{12.23}$$

From the consistency condition (12.21) we have

$$\begin{aligned} \tau_n &= \frac{y(x_{n-1}) - y_n(x_{n-1})}{h} - f(x_{n-1}, y(x_{n-1})) \\ &\quad + (\Phi(x_{n-1}, y(x_{n-1}); 0) - \Phi(x_{n-1}, y(x_{n-1}); h)). \end{aligned} \tag{12.24}$$

According to the Mean Value Theorem, Theorem A.3, the expression in the first bracket is equal to $y'(x) - y_n'(x)$, where $x \in [x_{n-1}, x_n]$. By Picard’s Theorem, y' is continuous on the closed interval $[x_0, X_0]$; therefore, it is uniformly continuous on this interval. Hence, for each $\epsilon > 0$ there exists $h(\epsilon)$ such that

$$|y'(x) - y_n'(x)| < \epsilon \quad \text{for } h < h(\epsilon), \quad n = 0, 1, \dots, N - 1.$$

Also, since $\Phi(\cdot, \cdot; \cdot)$ is a continuous function on the closed set $D = [0, h]$

for any pair of points $(x_0, y(x_0))$, $(x_1, y(x_1))$ on the solution curve.

12.4 An implicit one-step method

A one-step method with second-order accuracy is the **trapezium rule method**

$$y_1 = y_0 + \frac{h}{2} [f(x_0, y_0) + f(x_1, y_1)]. \tag{12.26}$$

This method is easily motivated by writing

$$y(x_1) - y(x_0) = \int_{x_0}^{x_1} y'(x) dx,$$

and approximating the integral by the trapezium rule. Since the right-hand side involves the integral of the function $y'(x) = f(x, y(x))$ we see at once from (7.6) that the truncation error

$$T = \frac{y(x_1) - y(x_0)}{h} - [f(x_0, y(x_0)) + f(x_1, y(x_1))]$$

of the trapezium rule method satisfies the bound

$$|T| \leq \frac{1}{12} h^2 M, \quad \text{where } M = \max_{x \in [x_0, x_1]} |y''(x)|. \tag{12.27}$$

The important difference between this method and Euler's method is that the value y_1 appears on both sides of (12.26). To calculate y_1 from the known y_0 therefore requires the solution of an equation, which will usually be nonlinear. This additional complication means an increase in the amount of computation required, but not usually a very large increase. The equation (12.26) is easily solved for y_1 by Newton's method, assuming that the derivative f'/y can be calculated quickly; as a starting point for the Newton iteration the obvious estimate

$$y_1 + hf(x_0, y_0),$$

will usually be close, and a couple of iterations will then suffice.

Methods of this type, which require the solution of an equation to determine the new value y_1 , are known as **implicit methods**.

Writing the trapezium rule method in the standard form (12.13) we see that

$$\begin{aligned} h\Phi(x_1, y_1; h) &= \frac{h}{2} [f(x_0, y_0) + f(x_1, y_1)] \\ &= \frac{h}{2} [f(x_0, y_0) + f(x_1, y_0 + h\Phi(x_1, y_1; h))]. \end{aligned} \tag{12.28}$$

Hence, the function Φ is also defined in an implicit form.

In order to employ Theorem 12.2 to estimate the error in the trapezium rule method we need a value for the Lipschitz constant $L_{\%}$. From (12.28) we find that

$$\Phi(x, u; h) - \Phi(x, v; h) = -f(x, u) - f(x + h, u + h\Phi(x, u; h)) \\ + f(x, v) + f(x + h, v + h\Phi(x, v; h)).$$

Hence,

$$\begin{aligned} & \Phi(x, u; h) - \Phi(x, v; h) \\ & - f(x, u) - f(x, v) \\ & + f(x + h, u + h\Phi(x, u; h)) - f(x + h, v + h\Phi(x, v; h)) \\ & -L |u - v| \\ & +L |u + h\Phi(x, u; h) - v - h\Phi(x, v; h)| \\ & -L |u - v| + -L |u - v| + -L |h\Phi(x, u; h) - h\Phi(x, v; h)|. \end{aligned}$$

This shows that

$$| \Phi(x, u; h) - \Phi(x, v; h) | \leq L |u - v|,$$

and, therefore,

$$L_{\%} = \frac{L}{1 - hL}, \quad \text{provided that } hL < 1.$$

Consequently, (12.16) and (12.27) imply that the global error in the trapezium rule method is $O(h^2)$, as h tends to 0.

Figure 12.3 depicts the results of some numerical calculations on the interval $x \in [0, 1.6]$ for the same problem as in Figure 12.2. The step sizes are 0.4 and 0.2, larger than for Euler’s method; nevertheless we see a much reduced error in comparison with Euler’s method, and also how the reduction in the step size h by a factor of 2 gives a reduction in the error by a factor of about 4, as predicted by our error analysis.

12.5 Runge–Kutta methods

Euler’s method is only first-order accurate; nevertheless, it is simple and cheap to implement because, to obtain y_{n+1} from y_n , we only require a single evaluation of the function f , at (x_n, y_n) . Runge–Kutta methods aim to achieve higher accuracy by sacrificing the efficiency of Euler’s method through re-evaluating $f(x, y)$ at points intermediate between



$$\begin{aligned}
 y(x) &= f + f y = f + f f, \\
 y(x) &= f + f f + (f + f f)f + f(f + f f),
 \end{aligned}$$

and so on; in these expressions the subscripts x and y denote partial derivatives, and all functions appearing on the right-hand sides are to be evaluated at $(x, y(x))$. We also need to expand $\Phi(x, y(x); h)$ in powers of h , giving (with the same notational conventions as before)

$$\begin{aligned}
 \Phi(x, y(x); h) &= a f + b f + h f + h f f + \frac{1}{2} (h)^2 f \\
 &\quad + \frac{1}{2} h^2 f f + \frac{1}{6} (h)^3 f f f + \frac{1}{24} (h)^4 f f f f.
 \end{aligned}$$

Thus, we obtain the truncation error in the form

$$\begin{aligned}
 T &= \frac{y(x+h) - y(x)}{h} - \Phi(x, y(x); h) \\
 &= f + \frac{1}{2} h (f + f f) \\
 &\quad + \frac{1}{6} h^2 [f + 2f f + f f + f (f + f f)] \\
 &\quad + \frac{1}{24} h^3 [a f + b (f + h f + h f f + \frac{1}{2} (h)^2 f \\
 &\quad + h f f + \frac{1}{6} (h)^3 f f f) + \frac{1}{24} (h)^4 f f f f].
 \end{aligned}$$

As $1 - a - b = 0$, the term $(1 - a - b)f$ is equal to 0. The coefficient of the term in h is

$$\frac{1}{2} (f + f f) - b f - b f f$$

which vanishes for all functions f provided that

$$b = \frac{1}{2} (1 + b).$$

The method is therefore second-order accurate if

$$a = 1 - \frac{1}{2} b, \quad b = \frac{1}{2} (1 + b), \quad c = 0,$$

showing that there is a one-parameter family of second-order methods of this form, parametrised by $c = 0$. The truncation error of the method then becomes

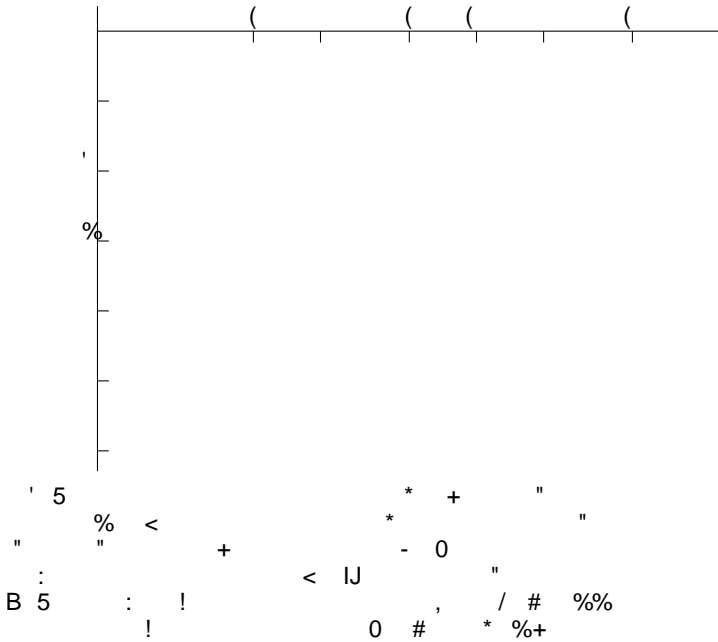
$$\begin{aligned}
 T &= \frac{1}{6} h^2 (-\frac{1}{2} b)(f + f f) + (-\frac{1}{24} b^2) f f \\
 &\quad + \frac{1}{24} h^3 (f f + f f f) + \frac{1}{24} (h)^4 f f f f.
 \end{aligned} \tag{12.32}$$

Evidently there is no choice of the free parameter b which will make this method third-order accurate for all functions f ; this can be seen, for example, by considering the initial value problem $y' = y, y(0) = 1$, and noting that in this case (12.32), with $f(x, y) = y$, yields

$$T = -\frac{1}{6} h^2 y(x) + \frac{1}{24} (h)^3 = -\frac{1}{6} h^2 e^x + \frac{1}{24} (h)^3.$$

Two examples of second-order Runge–Kutta methods of the form (12.29)–(12.31) are the modified Euler method and the improved Euler method.

(a)



To illustrate the behaviour of the one-step methods which we have discussed, Figure 12.4 shows the errors in the calculation of $y(1.6)$, where $y(x)$ is the solution to the problem (12.20) on the interval $[0, 1.6]$. The horizontal axis indicates N , the number of equally spaced mesh points used in the interval $(0, 1.6]$, on a logarithmic scale, and the vertical axis shows $\ln e = \ln |y(1.6) - y|$. The three methods employed are Euler's method, the trapezium rule method, and the classical Runge-Kutta method (12.33). The three lines show clearly the improved accuracy of the higher-order methods, and the rate at which the accuracy improves as N increases.

12.6 Linear multistep methods

While Runge-Kutta methods give an improvement over Euler's method in terms of accuracy, this is achieved by investing additional computational effort; in fact, Runge-Kutta methods require more evaluations of $f(x, y)$ than would seem necessary. For example, the fourth-order method involves four function evaluations per step. For comparison, by considering three consecutive points x_{n-2} , $x_{n-1} = x_{n-2} + h$, $x_n = x_{n-2} + 2h$, integrating the differential equation between x_{n-2} and x_{n-1} ,

yields

$$y(x_{n+1}) = y(x_n) + \int_{x_n}^{x_{n+1}} f(x, y(x)) dx,$$

and applying Simpson's rule to approximate the integral on the right-hand side then leads to the method

$$y_{n+1} = y_{n-1} + \frac{1}{3}h [f(x_{n-1}, y_{n-1}) + 4f(x_n, y_n) + f(x_{n+1}, y_{n+1})], \quad (12.34)$$

requiring only three function evaluations per step. In contrast with the one-step methods considered in the previous section where only a single value y_n was required to compute the next approximation y_{n+1} , here we need *two* preceding values, y_n and y_{n-1} , to be able to calculate y_{n+1} , and therefore (12.34) is *not* a one-step method.

In this section we consider a class of methods of the type (12.34) for the numerical solution of the initial value problem (12.1), (12.2), called **linear multistep methods**.

Given a sequence of equally spaced mesh points (x_n) with step size h , we consider the general **linear k-step method**

$$y_{n+1} = h \sum_{j=0}^k \alpha_j f(x_n, y_n), \quad (12.35)$$

where the coefficients $\alpha_0, \dots, \alpha_k$ and β_0, \dots, β_k are real constants. In order to avoid degenerate cases, we shall assume that $\beta_k \neq 0$ and that α_0 and β_0 are not both equal to 0. If $\beta_0 = 0$, then y_{n+1} is obtained explicitly from previous values of y_n and $f(x_n, y_n)$, and the k -step method is then said to be **explicit**. On the other hand, if $\beta_0 \neq 0$, then y_{n+1} appears not only on the left-hand side but also on the right, within $f(x_n, y_n)$; due to this implicit dependence on y_{n+1} the method is then called **implicit**. The method (12.35) is called *linear* because it involves only linear combinations of the y_n and the $f(x_n, y_n)$, $j = 0, 1, \dots, k$; for the sake of notational simplicity, henceforth we shall often write f instead of $f(x_n, y_n)$.

Example 12.4 *We have already seen an example of a linear two-step method in (12.34); here we present further examples of linear multistep methods.*

(a) Euler's method is a trivial case: it is an explicit linear one-step

method. The **implicit Euler method**

$$y_{n+1} = y_n + hf(x_{n+1}, y_{n+1}) \tag{12.36}$$

is an implicit linear one-step method. Another trivial example is the **trapezium rule method**, given by

$$y_{n+1} = y_n + \frac{1}{2}h(f_n + f_{n+1});$$

it, too, is an implicit linear one-step method.

(b) The **Adams–Bashforth method**

$$y_{n+1} = y_n + \frac{1}{24}h(55f_n - 59f_{n-1} + 37f_{n-2} - 9f_{n-3})$$

is an example of an explicit linear four-step method, while the **Adams–Moulton method**

$$y_{n+1} = y_n + \frac{1}{24}h(9f_{n+1} + 19f_n - 5f_{n-1} + 9f_{n-2})$$

is an implicit linear three-step method.

There are systematic ways of generating linear multistep methods, but these constructions will not be discussed here. Instead, we turn our attention to the analysis of linear multistep methods and introduce the concepts of (*zero-*) *stability*, *consistency* and *convergence*. The significance of these properties cannot be overemphasised: the failure of any of the three will render the linear multistep method practically useless.

12.7 Zero-stability

As is clear from (12.35) we need k starting values, y_0, \dots, y_{k-1} , before we can apply a linear k -step method to the initial value problem (12.1), (12.2): of these, y_0 is given by the initial condition (12.2), but the others,

1 K , % B> K #.#@ : C # K 6 #.@ \$
 4) #. # D , K || 6 \$
 1 1 5 6 5 ,, 1
 , 1 , 5 4 9 6 D , 1
 H 5 6! % 5 , 6!
 2 94 1 2 , O) 0 8 .%
 , 8 2
 O) \$ - 5 6 #..*4
 3 94 4' 2 (+ ' .% 5 6 1 ,
 #@ ?4

y_1, \dots, y_n , have to be computed by other means: say, by using a

Before stating the main theorem of this section, we recall a classical result from the theory of k th-order linear recurrence relations.

Lemma 12.1 *Consider the k th-order homogeneous linear recurrence relation*

$$y_{n+k} + a_{k-1}y_{n+k-1} + \dots + a_1y_{n+1} + a_0y_n = 0, \quad n = 0, 1, 2, \dots, \quad (12.38)$$

with $y_0 = 0, \dots, y_{k-1} = 0, \dots$, $j = 0, 1, \dots, k$, and the corresponding characteristic polynomial

$$p(z) = z^k + a_{k-1}z^{k-1} + \dots + a_1z + a_0.$$

Let z_1, \dots, z_r , \dots, z_k , be the distinct roots of the polynomial $p(z)$, and let $m_1, \dots, m_r, \dots, m_k$ denote the multiplicity of z_j , with $m_1 + \dots + m_r = k$. If a sequence (y_n) of complex numbers satisfies (12.38), then

$$y_n = p_1(n)z_1^n + \dots + p_r(n)z_r^n, \quad \text{for all } n \geq 0, \quad (12.39)$$

where $p_j(\cdot)$ is a polynomial in n of degree $m_j - 1, 1 \leq j \leq r$. In particular, if all roots are simple, that is $m_j = 1, 1 \leq j \leq k$, then the $p_j, r = 1, \dots, k$, are constants.

Proof We give a sketch of the proof. Let us first consider the case when all of the (distinct) roots z_1, z_2, \dots, z_k are simple. As, by assumption, $a_0 \neq 0$, none of the roots is equal to 0. It is then easy to verify by direct substitution that, since $p(z_r) = 0, r = 1, 2, \dots, k$, each of the sequences $(y_n) = (z_r^n), r = 1, 2, \dots, k$, satisfies (12.38).

In order to prove that any solution (y_n) of (12.38) can be expressed as a linear combination of the sequences $(z_1^n), (z_2^n), \dots, (z_k^n)$, it suffices to show that these k sequences are linearly independent. To do so, let us suppose that

$$C_1z_1^n + C_2z_2^n + \dots + C_kz_k^n = 0, \quad \text{for all } n = 0, 1, 2, \dots$$

Then, in particular,

$$\begin{aligned} C_1 + C_2 + \dots + C_k &= 0, \\ C_1z_1 + C_2z_2 + \dots + C_kz_k &= 0, \\ &\dots\dots\dots \\ C_1z_1^{k-1} + C_2z_2^{k-1} + \dots + C_kz_k^{k-1} &= 0. \end{aligned}$$

¹ 9 :) # 1 8 / 4 #*C # 1 4 0 ,) 5 ' 6 ! " #@? 4

The matrix of this system of k simultaneous linear equations for the k unknowns C_0, C_1, \dots, C_{k-1} has the determinant

$$\Delta = \begin{vmatrix} 1 & 1 & \dots & 1 \\ z_0 & z_1 & \dots & z_{k-1} \\ \dots & \dots & \dots & \dots \\ z_0^{k-1} & z_1^{k-1} & \dots & z_{k-1}^{k-1} \end{vmatrix},$$

known as the Vandermonde determinant, and $\Delta = \prod_{0 \leq i < j < k} (z_j - z_i)$. Since the roots are distinct, $\Delta \neq 0$, so the matrix of the system is nonsingular. Therefore $C_0 = C_1 = \dots = C_{k-1} = 0$ is the unique solution, which then means that the sequences $(z_0^n), (z_1^n), \dots, (z_{k-1}^n)$ are linearly independent.

Now, suppose that (y^n) is any solution of (12.38); as $\Delta \neq 0$, there exists a unique set of k constants, C_0, C_1, \dots, C_{k-1} , such that

$$y^n = C_0 z_0^n + C_1 z_1^n + \dots + C_{k-1} z_{k-1}^n, \quad n = 0, 1, \dots, k-1. \tag{12.40}$$

Substituting these equalities into (12.38) for $n = 0$, we conclude that

$$\begin{aligned} 0 &= y^0 + \dots - (C_0 z_0^0 + \dots + C_{k-1} z_{k-1}^0) + \\ &\quad + (C_0 z_0^1 + \dots + C_{k-1} z_{k-1}^1) \\ &= y^0 + C_0 (z_0^1 - z_0^0) + \dots + C_{k-1} (z_{k-1}^1 - z_{k-1}^0) \\ &= (y^0 - (C_0 z_0^0 + \dots + C_{k-1} z_{k-1}^0)). \end{aligned}$$

As $\Delta \neq 0$, it follows that

$$y^n = C_0 z_0^n + \dots + C_{k-1} z_{k-1}^n,$$

which, together with (12.40), proves (12.39) for $0 \leq n \leq k$ in the case of simple roots. Next, we select $n = 1$ in (12.38) and proceed in the same manner as in the case of $n = 0$ discussed above to show that (12.39) holds for $0 \leq n \leq k + 1$. Continuing in the same way, we deduce by induction that (12.39) holds for all $n \geq 0$.

In the case when (z) has repeated roots, the proof is similar, except that instead of (z^r) , $r = 1, 2, \dots, n$, the following k sequences are used:

$$\begin{aligned} &(z^r), \\ &(nz^r), \\ &\dots \\ &(n(n-1)\dots(n-m+2)z^r), \quad r = 1, 2, \dots, k. \end{aligned} \tag{12.41}$$

These can be shown to satisfy (12.38) by direct substitution on noting that $(z^r) = (z^r) = \dots = z^{r-1} (z) = 0$, given that z is a root of

(z) of multiplicity m, r = 1, 2, ..., . The linear independence of the sequences (12.41) follows as before, except instead of $\prod_{j=1}^m (z - z_j)^{m_j}$, the value of the corresponding determinant is now

$$\Delta = \prod_{j=1}^m \prod_{k=0}^{m_j-1} (z - z_j)^{m_j - k} = \prod_{j=1}^m (z - z_j)^{m_j!}$$

where 0! = 1, m! = m!(m - 1)!...1! for m = 1, 2, As the roots z₁, z₂, ..., z_m are distinct, we have that Δ ≠ 0, and therefore the sequences (12.41) are linearly independent. The rest of the argument is identical as in the case of simple roots. □

Now, we are ready to state the main result of this section.

Theorem 12.4 (Root Condition) *A linear multistep method is zero-stable for any initial value problem of the form (12.1), (12.2), where f satisfies the hypotheses of Picard's Theorem, if, and only if, all roots of the first characteristic polynomial of the method are inside the closed unit disc in the complex plane, with any which lie on the unit circle being simple.*

The algebraic stability condition contained in this theorem, namely that *the roots of the first characteristic polynomial lie in the closed unit disc and those on the unit circle are simple*, is often called the **Root Condition**.

Proof of theorem
 y = 0:

Consider the method (12.35), applied to

$$y_{k+1} + \dots + y_k + y_{k-1} = 0. \tag{12.42}$$

According to Lemma 12.1, every solution of this kth-order linear recurrence relation has the form

$$y_k = \sum_{j=1}^m p_j(n) z_j^k, \tag{12.43}$$

where z_j is a root, of multiplicity m_j - 1, of the first characteristic polynomial of the method, and the polynomial p_j has degree m_j - 1, 1 ≤ j ≤ m. Clearly, if |z_τ| > 1 for some τ, then there are starting values y₀, y₁, ..., y_{k-1} for which the corresponding solution grows like

$$|y_k| \sim B |z_\tau|^k, \quad B > 0, \quad \tau = 1, \dots, m$$

$z = 1$, and if $z = 1$ and the multiplicity is $m > 1$, then there is a solution growing like n^{-m} . In either case there are solutions that grow unboundedly as $n \rightarrow \infty$, i.e., as $h \rightarrow 0$ with nh fixed. Considering starting values y_0, y_1, \dots, y_{m-1} which give rise to such an unbounded solution (y_n), and starting values $z_0 = z_1 = \dots = z_{m-1} = 0$ for which the corresponding solution of (12.42) is (z_n) with $z_n = 0$ for all n , we see that (12.37) cannot hold. To summarise, if the Root Condition is violated, then the method is not zero-stable.

The proof that the Root Condition is sufficient for zero-stability is long and technical, and will be omitted here. For details, the interested reader is referred to Theorem 3.1 on page 353 of W. Gautschi, *Numerical Analysis: an Introduction*, Birkhäuser, Boston, MA, 1997. □

Example 12.5 *We shall explore the zero-stability of the methods from Example 12.4 using the Root Condition.*

(a) The Euler method and the implicit Euler method have first characteristic polynomial $(z) = z - 1$ with simple root $z = 1$, so both methods are zero-stable. The same is true of the trapezium rule method.

(b) The Adams–Bashforth and Adams–Moulton methods considered in Example 12.4 have first characteristic polynomials, respectively, $(z) = z(z - 1)$ and $(z) = z(z - 1)$. These have multiple root $z = 0$ and simple root $z = 1$, and therefore both methods are zero-stable.

(c) The three-step method

$$11y_{n+3} + 27y_{n+2} - 27y_{n+1} - 11y_n = 3h(f_{n+3} + 9f_{n+2} + 9f_{n+1} + f_n) \tag{12.44}$$

is *not* zero-stable. Indeed, the corresponding first characteristic polynomial $(z) = 11z^3 + 27z^2 - 27z - 11$ has roots at $z = 1$, $z = 0.32$, $z = -3.14$, so $|z| > 1$.

(d) The first characteristic polynomial of the three-step method

$$y_{n+3} + y_{n+2} - y_{n+1} - y_n = 2h(f_{n+3} + f_{n+2})$$

is $(z) = z^3 + z^2 - z - 1 = (z + 1)(z - 1)$, which has roots $z = -1$, $z = 1$. The first of these is a double root lying on the unit circle; therefore, the method is *not* zero-stable.

12.8 Consistency

In this section we consider the accuracy of the linear k-step method (12.35). For this purpose, as in the case of one-step methods, we introduce the notion of truncation error. Thus, suppose that y is a solution to the ordinary differential equation (12.1). The truncation error of (12.35) is then defined as follows:

$$\tau = \frac{[y(x_{j+1}) - h f(x_j, y(x_j))]}{h}. \tag{12.45}$$

Of course, the definition requires implicitly that $y(x_{j+1}) - h f(x_j, y(x_j)) = 0$. Again, as in the case of one-step methods, the truncation error can be thought of as the residual that is obtained by inserting the solution of the differential equation into the formula (12.35) and scaling this residual appropriately (in this case dividing through by h), so that τ resembles $y - f(x, y(x))$.

Definition 12.4 *The numerical method (12.35) is said to be consistent with the differential equation (12.1) if the truncation error defined by (12.45) is such that for any $\epsilon > 0$ there exists an $h(\epsilon)$ for which*

$$|\tau| < \epsilon \quad \text{for } 0 < h < h(\epsilon),$$

and any $k + 1$ points $(x_0, y(x_0)), \dots, (x_k, y(x_k))$ on any solution curve in D of the initial value problem (12.1), (12.2).

Now, let us suppose that the solution to the differential equation is sufficiently smooth, and let us expand the expressions $y(x_{j+1})$ and $f(x_j, y(x_j)) = y'(x_j)$ into Taylor series about the point x_j . On substituting these expansions into the numerator in (12.45) we obtain

$$\tau = \frac{1}{h} [C_0 y(x_j) + C_1 h y'(x_j) + C_2 h^2 y''(x_j) + \dots - h f(x_j, y(x_j))] \tag{12.46}$$

where

$$\begin{aligned} C_0 &= y(x_j) - h f(x_j, y(x_j)), \\ C_1 &= y'(x_j) - f(x_j, y(x_j)), \\ C_2 &= \frac{1}{2} y''(x_j) - f'(x_j, y(x_j)), \\ &\vdots \\ C_j &= \frac{1}{j!} y^{(j)}(x_j) - \frac{d^j f}{dx^j}(x_j, y(x_j)). \end{aligned} \tag{12.47}$$

For consistency we need that, as $h \rightarrow 0$ and $n \rightarrow \infty$ with $x_n = x$ $[x_n, X_0]$, the truncation error T_n tends to 0. This requires that $C_0 = 0$ and $C_1 = 0$ in (12.46). In terms of the characteristic polynomials this consistency requirement can be restated in compact form as

$$P(z) = 0 \quad \text{and} \quad P'(z) = P'(1) = 0.$$

Let us observe that, according to this condition, if a linear multistep method is consistent, then it has a *simple* root on the unit circle at $z = 1$; thus, the Root Condition is not violated by this root.

Definition 12.5 *The numerical method (12.35) is said to have order of accuracy p , if p is the largest positive integer such that, for any sufficiently smooth solution curve in D of the initial value problem (12.1), (12.2), there exist constants K and h_0 such that*

$$T_n \leq Kh^{p+1} \quad \text{for } 0 < h \leq h_0,$$

for any $k + 1$ points $(x_0, y(x_0)), \dots, (x_k, y(x_k))$ on the solution curve.

Thus, we deduce from (12.46) that the method is of order of accuracy p if, and only if,

$$C_0 = C_1 = \dots = C_p = 0 \quad \text{and} \quad C_{p+1} \neq 0.$$

In this case,

$$T_n = \frac{C_{p+1}}{(p+1)!} h^{p+1} y^{(p+1)}(x) + O(h^{p+2}).$$

The number $C_{p+1} / (p+1)!$ is called the **error constant** of the method.

Example 12.6 *Let us determine all values of the real parameter b , $b \neq 0$, for which the linear multistep method*

$$y_{n+4} + (2b - 3)y_{n+3} - y_n = h(b(f_{n+4} + f_n))$$

is zero-stable. We shall show that there exists a value of b for which the order of the method is 4, and that if the method is zero-stable for some value of b , then its order cannot exceed 2.

According to the Root Condition, this linear multistep method is zero-stable if, and only if, all roots of its first characteristic polynomial

$$\rho(z) = z^2 + (2b - 3)z - 1$$

belong to the closed unit disc, and those on the unit circle are simple.

Clearly, $\rho(1) = 0$; upon dividing $\rho(z)$ by $z - 1$ we see that $\rho(z)$ can be written in the following factorised form:

$$\rho(z) = (z - 1)(z - \alpha), \quad \text{where} \quad \alpha = z^{-2}(1 - b)z + 1.$$

Thus, the method is zero-stable if, and only if, all roots of the polynomial

$\rho(z)$ belong to the closed unit disc, and those on the unit circle are simple and differ from 1. Suppose that the method is zero-stable. It then follows that $b = 0$ and $b = 2$, since these values of b correspond to double roots of $\rho(z)$ on the unit circle, respectively, $z = 1$ and $z = -1$. Further, since the product of the two roots of $\rho(z)$ is equal to 1, both have modulus less than or equal to 1, and neither of them is equal to -1 , it follows that they must both be strictly complex; hence the discriminant of the quadratic polynomial $\rho(z)$ must be negative. That is, $4(1 - b)^2 - 4 < 0$. In other words, $b \in (0, 2)$.

Conversely, suppose that $b \in (0, 2)$. Then, the roots of $\rho(z)$ are

$$z_1 = 1, \quad z_2 = 1 - b + i \sqrt{1 - (b - 1)^2}.$$

Since $z_1 = 1$, $z_2 = 1$ and $z_2 = \bar{z}_1$, all roots of $\rho(z)$ lie on the unit circle and they are simple. Hence the method is zero-stable. To summarise, the method is zero-stable if, and only if, $b \in (0, 2)$.

In order to analyse the order of accuracy of the method, we note that, upon Taylor series expansion, its truncation error can be written in the form

$$\begin{aligned} T_n &= \frac{1}{(1)} - 1 - \frac{b}{6} h^2 y''(x) + \frac{1}{4}(6 - b)h^3 y'''(x) \\ &\quad + \frac{1}{120}(150 - 23b)h^4 y^{(4)}(x) + O(h^5), \end{aligned}$$

where $\rho(1) = 2b = 0$. If $b = 6$, then $T_n = O(h^4)$ and so the method is of order 4. As $b = 6$ does not belong to the interval $(0, 2)$, we deduce that the method is *not* zero-stable for $b = 6$.

Since zero-stability requires $b \in (0, 2)$, in which case $1 - b = 0$, it follows that if the method is zero-stable, then $T_n = O(h^4)$.

12.9 Dahlquist's theorems

An important result connecting the concepts of zero-stability, consistency and convergence of a linear multistep method was proved by the Swedish mathematician Germund Dahlquist.

Theorem 12.5 (Dahlquist's Equivalence Theorem) *For a linear k -step method that is consistent with the ordinary differential equation (12.1) where f is assumed to satisfy a Lipschitz condition, and with consistent starting values, zero-stability is necessary and sufficient for convergence. Moreover if the solution y has continuous derivative of order $p + 1$ and truncation error $\tau_n(h)$, then the global error of the method, $e_n = y(x_n) - y_n$, is also $O(h^p)$.*

The proof of this result is long and technical; for details of the argument, see Theorem 6.3.4 on page 357 of W. Gautschi, *Numerical Analysis: an Introduction*, Birkhäuser, Boston, MA, 1997, or Theorem 5.10 on page 244 of P. Henrici, *Discrete Variable Methods in Ordinary Differential Equations*, Wiley, New York, 1962.

By virtue of Dahlquist's theorem, if a linear multistep method is not zero-stable its global error cannot be made arbitrarily small by taking the mesh size h sufficiently small for any sufficiently accurate initial data. In fact, if the Root Condition is violated, then there exists a solution to the linear multistep method which will grow by an arbitrarily large factor in a fixed interval of x , however accurate the starting conditions are. This result highlights the importance of the concept of zero-stability and indicates its relevance in practical computations.

A second theorem by Dahlquist imposes a restriction on the order of accuracy of a zero-stable linear multistep method.

Theorem 12.6 (Dahlquist's Barrier Theorem) *The order of accuracy of a zero-stable k -step method cannot exceed $k + 1$ if k is odd, or $k + 2$ if k is even.*

A proof of this result will be found in Section 4.2 of Gautschi's book or in Section 5.2-8 of Henrici's book, cited above.

Theorem 12.6 makes it very difficult to choose a 'best' multistep method of a given order. Suppose, for example, that we consider five-step methods. The general five-step method involves 12 parameters, of

¹ 8, 5, 0, \$ 4 " A # # BSD % A & # , , 5

which 11 are independent: the method is obviously unaffected by multiplying all the parameters by a nonzero constant. Now it would be possible to construct a five-step method of order 10, by solving the 11 equations of the form $C = 0$, $q = 0, 1, \dots, 10$, where C is given in (12.47). But the Barrier Theorem states that this method would not be zero-stable, and the order of a zero-stable five-step method cannot exceed 6. There is a family of stable five-step methods of order 6, involving 4 free parameters, and there is no obvious way of deciding whether any one of these methods is better than the others.

Example 12.7 (i) *The Barrier Theorem says that when $k = 1$ the order of accuracy of a zero-stable method cannot exceed 2. The trapezium rule method has order 2, and is zero-stable.*

(ii) *The two-step method*

$$y_{n+2} - y_n = h(-f_n + -f_{n+1} + -f_{n+2})$$

is zero-stable, as the roots of the first characteristic polynomial, $(z) = z^2 - 1$, are 1 and -1 . A simple calculation shows that its order of accuracy is 4; by the Barrier Theorem, this is the highest order which could be achieved by a two-step method.

(iii) *The three-step method*

$$11y_{n+3} + 27y_{n+2} - 27y_{n+1} - 11y_n = 3h(f_n + 9f_{n+1} + 9f_{n+2} + f_{n+3})$$

has order 6. The Barrier Theorem therefore implies that this method is not zero-stable. We have already shown this in Example 12.5(c) using the Root Condition.

It is found that all the zero-stable k -step methods of highest possible order are *implicit*, with α_k nonzero.

12.10 Systems of equations

In this section we discuss the application of numerical methods to simultaneous systems of differential equations, which we shall write in the form

$$\frac{d}{dx} = (x, y).$$

Here \mathbf{y} is an m -component vector function of x , and \mathbf{f} is an m -component vector function of the independent variable x and the vector variable \mathbf{y} . In component form the system becomes

$$\frac{dy_j}{dx} = f_j(x, y_1, \dots, y_m), \quad j = 1, 2, \dots, m.$$

The system comprises m simultaneous differential equations. To single out a unique solution we need m side conditions, and we shall suppose that all these conditions are given at the same value of x , and have the form

$$y_j(x_0) = y_{j0},$$

or, in component form,

$$y_j(x_0) = y_{j0}, \quad j = 1, 2, \dots, m,$$

where the values of y_{j0} are given. This is called an initial value problem for a system of ordinary differential equations; we may also require a solution of the system on an interval $[a, b]$, with r conditions given at one end of the interval and $m - r$ conditions at the other end. This constitutes a boundary value problem, and requires different numerical methods which are considered in the next chapter.

All the numerical methods which we have discussed apply without change to systems of differential equations; it is only necessary to realise that we are dealing with vectors. For example, the first stage of the classical Runge-Kutta method (12.33) becomes

$$k_1 = \mathbf{f}(x_0, \mathbf{y}_0);$$

we must evaluate all the elements of the vector k_1 before proceeding to the next stage to calculate k_2 , and so on.

The most important difference which arises in dealing with a system of differential equations is in the practical use of an *implicit* multi-step method. As we have seen, this almost always requires an iterative method for the solution of an equation to determine y_{n+1} . Applying such a method to a system of differential equations now involves the solution of a system of equations, which will usually be nonlinear, to determine the elements of the vector \mathbf{y}_{n+1} . In real-life problems it is quite common to deal with systems of several hundred differential equations, and it then becomes very important to be sure that the improved efficiency of the implicit method justifies the very considerable extra work in each step of the process.

We shall not discuss the extension of our earlier analysis to deal with

systems of differential equations; in almost all cases we simply need to introduce vector notation, and replace the absolute value of a number by the norm of a vector. For example, in the proof of Theorem 12.2, (12.17) becomes

$$\|y_{n+1} - y_{n+1}^h\| \leq \|T_n\| \|y_n - y_n^h\|, \quad n = 0, 1, \dots, N-1,$$

where $\|\cdot\|$ is any norm on \mathbb{R}^m , with obvious definitions of the global error $\|y - y^h\|$ and the truncation error $\|T_n\|$. Similarly, Picard's Theorem and its proof, discussed at the beginning of the chapter in the case of a single ordinary differential equation, can be easily extended to an m -component system of differential equations by replacing the absolute value sign with a vector norm on \mathbb{R}^m throughout.

12.11 Stiff systems

The phenomenon of stiffness usually appears only in a system of differential equations, but we begin by discussing an almost trivial example of a single equation,

$$y' = -\lambda y, \quad y(0) = y_0,$$

where λ is a constant. The solution of this equation is evidently $y(x) = y_0 \exp(-\lambda x)$. When $\lambda < 0$ the absolute value of the solution is exponentially decreasing, so it is sensible to require that the absolute value of our numerical solution also decreases. It is very easy to give expressions for the result of a numerical solution using Euler's method and the implicit Euler method (12.36). They are, respectively,

$$y_n^E = (1 - h\lambda)^{-n} y_0, \quad y_n^I = (1 - h\lambda)^{-n} y_0.$$

When $\lambda < 0$ and $h > 0$, we have $(1 - h\lambda) > 1$; therefore, the sequence (y_n^E) decreases monotonically with increasing n . On the other hand, for $\lambda < 0$ and $h > 0$,

$$1 + h\lambda < 1 \quad \text{if, and only if,} \quad 0 < h < 2/\lambda.$$

This gives the restriction $h < 2/\lambda$ on the size of h for which the sequence (y_n^E) decreases monotonically; if h exceeds $2/\lambda$, the numerical solution obtained by Euler's method will oscillate with increasing magnitude with increasing n and fixed $h > 0$, instead of converging to zero as $n \rightarrow \infty$.

We now consider the same two methods applied to the initial value problem for a system of differential equations of the form

$$y' = Ay, \quad y(0) = y_0,$$

where A

We consider the system where A is the 2×2 matrix

$$A = \begin{pmatrix} 8003 & 1999 \\ 23988 & 6004 \end{pmatrix}$$

and the initial condition is

$$y(0) = \begin{pmatrix} 1 \\ 4 \end{pmatrix}$$

The eigenvalues of A are $\lambda_1 = -7$ and $\lambda_2 = -14000$; the solution of the problem is

$$y(x) = \begin{pmatrix} e^{-7x} \\ 4e^{-14000x} \end{pmatrix}$$

Clearly, $\lim_{x \rightarrow \infty} y(x) = 0$.

The numerical solution uses 12 steps of size $h = 0.004$; the results are shown in Table 12.1. The second column gives the first component of the solution, $y_1(x) = e^{-7x}$, the third column shows the result from the implicit Euler method, and the last gives the result of the standard Euler method. The last column is a dramatic example of what happens when the step size h is too large; in this case $h = 56$. The numerical values given by the implicit Euler method have an error of a few units in the third decimal digit; to get the same accuracy from the Euler method would require a step size about 30 times smaller, and about 30 times as much work.

It is clear that the difficulty in the numerical example is caused by the size of the eigenvalue -14000 , but what is important is its size relative to the other eigenvalue. The special constant-coefficient system $y' = Ay$ is said to be **stiff** if all the eigenvalues of A have negative real parts, and if the ratio of the largest of the real parts to the smallest of the real parts is large. Most practical problems are nonlinear, and for such problems it is quite difficult to define precisely what is meant by stiffness. To begin with we may replace the system by a linearised approximation, the first terms of an expansion

$$y'(x) = f(x, y) + \frac{1}{x}(y, y'(x))(x, x) + J(x)(y(x), y'(x)) +$$

1) 5, 15, \$, S, 6, 1, J, G, \$
 J, E, 5, B, G, 5, D, 6, L, 1, 1
 (, 1, 5, E, 1, 24, 1, K4/4, : \$
 #@@# 6

Table 12.1. *The use of Euler’s method and the implicit Euler method to solve a stiff system.*

	$h = 10^{-7}$	$h = 10^{-8}$	$h = 10^{-9}$
$\ x(t) - x_{\text{exact}}(t)\ $	1.0×10^{-7}	1.0×10^{-8}	1.0×10^{-9}
$\ x(t) - x_{\text{Euler}}(t)\ $	1.0×10^{-7}	1.0×10^{-7}	1.0×10^{-7}
$\ x(t) - x_{\text{implicit Euler}}(t)\ $	1.0×10^{-7}	1.0×10^{-8}	1.0×10^{-9}

where J is the Jacobian matrix of the function f , whose (i, j) -entry is

$$(J(x))_{ij} = \frac{\partial f_i}{\partial x_j}(x, y).$$

We can then think of the system as being stiff if the eigenvalues of the matrix $J(x)$ have negative real parts and if the ratio of the largest of the real parts to the smallest is large. Although this gives some indication of the sort of problems which may cause difficulty, the behaviour of nonlinear systems is much more complicated than this. It is not difficult to construct examples in which all the eigenvalues of the Jacobian matrix have negative real parts, yet the norm of the solution of the differential equation is exponentially increasing as $x \rightarrow \infty$.

Even though any classification of nonlinear systems of differential equations into stiff and nonstiff, based only on monitoring the eigenvalues of $J(x)$, is somewhat simplistic, it does highlight some of the key difficulties. Stiff systems of differential equations arise in many application areas, a typical one being chemical engineering. For example, in parts of an oil refinery there may be a large number of substances undergoing chemical reactions with widely different reaction rates. These reaction rates correspond to the eigenvalues of the Jacobian matrix, and it is not unusual to find the ratio of the largest of the real parts to the smallest to be in excess of 10^6 . For such problems it is essential to find a numerical method which imposes no restriction on the step size;

Euler's method, which might require the restriction $10^{-10} h < 2$, would evidently be quite useless.

Application of the linear multistep method

$$y_{k+1} = h \sum_{j=0}^{k-1} \alpha_j f(x_k, y_k)$$

to the equation $y' = \lambda y$ leads to the k th-order linear recurrence relation

$$\sum_{j=0}^k \alpha_j y_{k-j} = 0. \tag{12.48}$$

The characteristic polynomial of the linear recurrence relation (12.48) is

$$p(z; h) = \sum_{j=0}^k \alpha_j z^j.$$

Alternatively, we can write this in terms of the first and second characteristic polynomials of the linear multistep method as

$$p(z; h) = \rho(z) - h \sigma(z).$$

In the present context, the polynomial $p(z; h)$ is usually referred to as the **stability polynomial** of the linear multistep method. According to Lemma 12.1, the general solution of the recurrence relation (12.48) can be expressed in terms of the distinct roots z_1, \dots, z_k , of $p(z; h)$. Letting m_j denote the multiplicity of the root z_j , $1 \leq j \leq k$, $m_1 + \dots + m_k = k$, we have that

$$y_k = \sum_{j=1}^k p_j(n) z_j^k, \tag{12.49}$$

where the polynomial $p_j(n)$ has degree $m_j - 1$, $1 \leq j \leq k$.

Clearly, the roots z_j are functions of h . For $h > 0$, with $\text{Re}(z_j) < 0$, the solution of the model problem

$$y' = \lambda y, \quad y(0) = y_0,$$

converges in y to 0 as $x \rightarrow \infty$. Thus, we would like to ensure that, when a linear multistep method is applied to this problem, the step size h can be chosen so that the resulting sequence of numerical approximations (y_k) exhibits an analogous behaviour as $n \rightarrow \infty$, that is, $\lim_{n \rightarrow \infty} y_k = 0$. By virtue of (12.49), this can be guaranteed by demanding that each root $z_j = z_j(h)$ has modulus less than 1.

Definition 12.6 A linear multistep method is said to be **absolutely stable** for a given value of h if each root $z = z(\lambda, h)$ of the associated stability polynomial $\rho(z; \lambda, h)$ satisfies $|z(\lambda, h)| < 1$.

Our aim is, therefore, to single out those values of h for which the linear multistep method is absolutely stable.

Definition 12.7 The **region of absolute stability** of a linear multistep method is the set of all points h in the complex plane for which the method is absolutely stable.

Ideally, the region of absolute stability of the method should admit all values of λ , $\text{Re}(\lambda) < 0$, so as to ensure that there is no limitation on the size of h , however large λ may be. This leads us to the next definition.

Definition 12.8 A linear multistep method is said to be **A-stable** if its region of absolute stability contains the negative (left) complex half-plane.

Unfortunately, the condition of A-stability is extremely demanding. Dahlquist has shown the following results which are collectively known as his **Second Barrier Theorem**:

- (i) No *explicit* linear multistep method is A-stable;
- (ii) No A-stable linear multistep method can have order greater than 2.
- (iii) The second-order A-stable linear multistep method with the smallest error constant is the trapezium rule method.

The trapezium rule method is a one-step method, so the associated stability polynomial has only one root, given by

$$z = \frac{1 - h\lambda}{1 + h\lambda}.$$

Evidently $|z| < 1$ if $\text{Re}(h\lambda) = h \text{Re}(\lambda) < 0$, so the trapezium rule method is indeed A-stable.

To construct useful methods of higher order we need to relax the condition of A-stability by requiring that the region of absolute stability should include a large part of the negative half-plane, and certainly that it contains the whole of the negative real axis.

The most efficient methods of this kind in current use are the **Backward Differentiation Formulae**, or BDF methods. These are the linear multistep methods (12.35) in which $\alpha_0 = 0, \alpha_j = 0 \quad j = 1, \dots, k-1, \alpha_k = 1, \alpha_{k+1} = 0$, and $\beta_j = 0$. Thus,

$$y_{k+1} + \alpha_1 y_k + \dots + \alpha_k y_1 = h f_{k+1}.$$

The coefficients are obtained by requiring that the order of accuracy of the method is as high as possible, *i.e.*, by making the coefficients C_j zero in (12.47) for $j = 0, 1, \dots, k$. For $k = 1$ this yields the implicit Euler method (BDF1), whose order of accuracy is, of course, 1; the method is A-stable. The choice of $k = 6$ results in the sixth-order, six-step BDF method (BDF6):

$$147y_{k+1} - 360y_k + 450y_{k-1} - 400y_{k-2} + 225y_{k-3} - 72y_{k-4} + 10y_{k-5} = 60hf_{k+1}. \quad (12.50)$$

Although the method (12.50) is not A-stable, its region of absolute stability includes the whole of the negative real axis (see Figure 12.5). For the intermediate values, $k = 2, 3, 4, 5$, we have the following k th-order, k -step BDF methods, respectively:

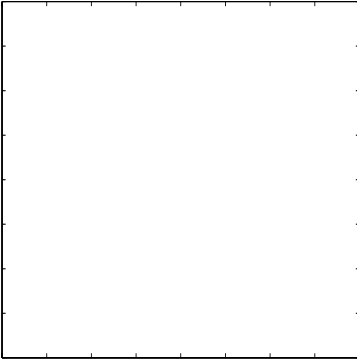
$$\begin{aligned} 3y_{k+1} - 4y_k + y_{k-1} &= 2hf_{k+1}, \\ 11y_{k+1} - 18y_k + 9y_{k-1} - 2y_{k-2} &= 6hf_{k+1}, \\ 25y_{k+1} - 48y_k + 36y_{k-1} - 16y_{k-2} + 3y_{k-3} &= 12hf_{k+1}, \\ 137y_{k+1} - 300y_k + 300y_{k-1} - 200y_{k-2} + 75y_{k-3} - 12y_{k-4} &= 60hf_{k+1}, \end{aligned}$$

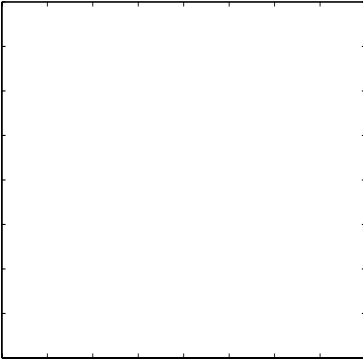
referred to as BDF2, BDF3, BDF4 and BDF5. Their regions of absolute stability are also shown in Figure 12.5. In each case the region of absolute stability includes the negative real axis. Higher-order methods of this type cannot be used, as all BDF methods, with $k > 6$, are zero-unstable.

12.12 Implicit Runge–Kutta methods

For Runge–Kutta methods absolute stability is defined in much the same way as for linear multistep methods; *i.e.*, by applying the method in question to the model problem $y' = \lambda y, y(0) = y_0, \lambda \in \mathbb{C}, \text{Re}(\lambda) < 0$, and demanding that the resulting sequence (y_n) converges to 0 as $n \rightarrow \infty$, with h held fixed. The set of all values of $h\lambda$ in the complex plane for which the method is absolutely stable is called the region of absolute stability of the Runge–Kutta method.

Classical Runge–Kutta methods are explicit, and are unsuitable for





where

$$k_i = f(x_i + hc_i, y_i + h a_i k_i), \quad 1 \leq i \leq s. \tag{12.51}$$

It is convenient to display the coefficients in a **Butcher tableau**

c	a	...	a
...
c	a	...	a
	b	...	b

The method is then defined by the matrix $A = (a_{ij})_{i,j=1}^s$, of order s , and the two vectors $c = (c_1, \dots, c_s)$ and $b = (b_1, \dots, b_s)$. For example, the classical four-stage Runge-Kutta method is defined by the tableau

0				
-	-			
-	0	-		
1	0	0	1	
	-	-	-	-

The 4×4 array representing the matrix A for this method, displayed in the upper right quadrant of the tableau, follows the usual notational convention that zero elements after the last nonzero element in each row of the matrix A are omitted.

This is an explicit method, shown by the fact that the matrix A is *strictly lower triangular*, with $a_{ij} = 0$ when $1 \leq i \leq j \leq 4$. Each value k_i can therefore be calculated in sequence, all the quantities on the right-hand side of (12.51) being known.

It is not difficult to construct s -stage implicit methods which are A -stable. For example, this can be done by choosing the coefficients c_i and b_i to be the quadrature points and weights respectively in the Gauss quadrature formula for the evaluation of

$$\int_a^b g(x) dx \approx \sum_{i=1}^s b_i g(c_i).$$

The numbers a_i can then be chosen so that the method has order $2s$, and is A -stable.

For example, the array

$$\begin{array}{ccc|ccc}
 -(3 & & 3) & & - & & -(3 & 2 & 3) \\
 -(3 + & & 3) & -(3 + 2 & 3) & & - & & \\
 \hline
 & & & - & & & - & &
 \end{array}$$

defines a 2-stage A-stable method of order 4.

However, there is a heavy price to pay for using implicit methods of this kind, as we now have to calculate all the numbers k_i , $i = 1, 2, \dots, s$, simultaneously, not in succession. For a system of m differential equations an implicit linear multistep method requires the solution of m simultaneous equations at each step; an s -stage implicit Runge–Kutta method requires the solution of sm simultaneous equations. This is a considerable increase in cost, and the general implicit Runge–Kutta methods cannot compete in efficiency with the Backward Differentiation Formulae such as (12.50); their use is almost exclusively limited to stiff systems of ODEs.

The overall computational effort can be somewhat reduced by using **diagonally implicit Runge–Kutta** (or DIRK) methods, in which the matrix A is lower triangular, so that $a_{ij} = 0$ if $j > i$. A further improvement in efficiency is possible by requiring in addition that all the diagonal elements a_{ii} are the same; unfortunately it has proved difficult to construct such methods with order greater than 4.

12.13 Notes

In this chapter we have only been able to introduce some of the basic ideas in what has become a vast area of numerical analysis. In particular we have not discussed the practical implementation of the various methods. The questions of how to choose the step size h to obtain efficiently a prescribed accuracy, and when and how to adjust h during the course of the calculation, are dealt with in the following books.

/ " & \$ 3 " , *Solving Ordinary Differential Equations I: Nonstiff Problems*, Second Edition, Springer Series in Computational Mathematics, 8, Springer, Berlin, 1993.
) # , *A First Course in the Numerical Analysis of Differential Equations*, Cambridge University Press, Cambridge, 1996.
 , *Numerical Methods for Ordinary Differential Systems*, John Wiley & Sons, Chichester, 1991.

For a study of dynamical systems and their numerical analysis, with focus on long-time behaviour, we refer to

) + &) % , *Dynamical Systems and Numerical Analysis*, Cambridge University Press, Cambridge, 1999.

The numerical solution of stiff initial value problems for systems of ordinary differential equations is discussed in

 / , *Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems*, Springer Series in Computational Mathematics, 14, Springer, Berlin, 1991.

An extensive survey of the theory of Runge–Kutta and linear multistep methods is found in

 ! , *The Numerical Analysis of Ordinary Differential Equations. Runge–Kutta and General Linear Methods*, Wiley-Interscience, John Wiley & Sons, Chichester, 1987.

Satisfactory theoretical treatment of nonlinear systems of differential equations from the point of view of stiffness requires the development of a genuinely nonlinear stability theory which does not involve the rather dubious idea of defining stiffness through linearisation based on the ‘frozen Jacobian matrix’. We close by mentioning just one concept in this direction – that of *algebraic stability*. Given a Runge–Kutta method with Butcher tableau

$$\begin{array}{c|c} & A \\ \hline b & \end{array}$$

we define the matrices

$$B = \text{diag}(b_1, b_2, \dots, b_s) \quad \text{and} \quad M = BA + A B^{-1}.$$

The method is said to be **algebraically stable** if the matrices B and M are both positive semidefinite, *i.e.*, $B \succeq 0$ and $M \succeq 0$ for all

 . Algebraic stability can be seen to ensure that approximations to solutions of nonlinear systems of differential equations exhibit acceptable numerical behaviour. For example, the Gauss–Runge–Kutta methods discussed in the last section are algebraically stable. For further details, see, for example,

 , -- , *Stability of Runge–Kutta Methods for Stiff Nonlinear Differential Equations*, North-Holland, Amsterdam, 1984.

Exercises

12.1 Verify that the following functions satisfy a Lipschitz condition on the respective intervals and find the associated Lipschitz constants:

- (a) $f(x, y) = 2yx^{-2}$, $x \in [1, \infty)$;
- (b) $f(x, y) = e^{-x} \tan^{-1} y$, $x \in [1, \infty)$;
- (c) $f(x, y) = 2y(1 + y)^{-1} (1 + e^{-x})^{-1}$, $x \in (0, \infty)$.

12.2 Suppose that m is a fixed positive integer. Show that the initial value problem

$$y' = y^m, \quad y(0) = 0,$$

has infinitely many continuously differentiable solutions. Why does this not contradict Picard's Theorem?

12.3 Write down the solution y of the initial value problem

$$y' = py + q, \quad y(0) = 1,$$

where p and q are constants. Suppose that the method in the proof of Picard's Theorem is used to generate the sequence of approximations $y_n(x)$, $n = 0, 1, 2, \dots$; show that $y_n(x)$ is a polynomial of degree n , and consists of the first $n + 1$ terms in the series expansion of $y(x)$ in powers of x .

12.4 Show that Euler's method fails to approximate the solution $y(x) = (4x/5)^{1/4}$ of the initial value problem $y' = y^{-3/4}$, $y(0) = 0$. Justify your answer.

Consider approximating the same problem with the implicit Euler method. Show that there is a solution of the form $y = (C h)^{1/4}$, $n = 0, 1, 2, \dots$, with $C = 0$ and $C = 1$ and $C > 1$ for all $n \geq 2$.

12.5 Write down Euler's method for the solution of the problem

$$y' = xe^{-y} - 5y, \quad y(0) = 0$$

on the interval $[0, 1]$ with step size $h = 1/N$. Denoting by y_N the resulting approximation to $y(1)$, show that $y_N \rightarrow y(1)$ as $N \rightarrow \infty$.

12.6 Consider the initial value problem

$$y' = \ln \ln(4 + y), \quad x \in [0, 1], \quad y(0) = 1,$$

and the sequence (y_n) , $n = 0, 1, \dots, N - 1$, generated by the Euler method

$$y_{n+1} = y_n + h \ln \ln(4 + y_n), \quad n = 0, 1, \dots, N - 1, \quad y_0 = 1,$$

using the mesh points $x_n = nh$, $n = 0, 1, \dots, N$, with spacing $h = 1/N$.

(i) Let T_n denote the truncation error of Euler's method for this initial value problem at the point $x = x_n$. Show that $T_n \leq h/4$.

(ii) Verify that

$$y(x_{n+1}) - y(x_n) = (1 + hL)y(x_n) - y(x_n) + hT_n$$

for $n = 0, 1, \dots, N - 1$, where $L = 1/(2 \ln 4)$.

(iii) Find a positive integer N_0 , as small as possible, such that

$$\max_{n=0, \dots, N-1} |y(x_{n+1}) - y(x_n)| < 10^{-6}$$

whenever $N \geq N_0$.

12.7 Define the truncation error T of the trapezium rule method

$$y_{n+1} = y_n + \frac{1}{2}h(f(x_n) + f(x_{n+1}))$$

for the numerical solution of $y' = f(x, y)$ with $y(0) = y_0$ given, where $f = f(x, y)$ and $h = x_{n+1} - x_n$.

By integrating by parts the integral

$$\int_{x_n}^{x_{n+1}} (x - x_n)(x - x_{n+1})y''(x)dx,$$

or otherwise, show that

$$T = \frac{1}{12}h^3 y'''(\xi)$$

for some ξ in the interval (x_n, x_{n+1}) , where y is the solution of the initial value problem.

Suppose that f satisfies the Lipschitz condition

$$|f(x, u) - f(x, v)| \leq L|u - v|$$

for all real x, u, v , where L is a positive constant independent

of x , and that $y(x) \leq M$ for some positive constant M independent of x . Show that the global error $e_n = y(x_n) - y_n$ satisfies the inequality

$$e_n \leq e^{-\frac{1}{2}hL} \left(e^{-\frac{1}{2}hL} + e^{-\frac{1}{2}hL} \right) + \frac{1}{12} h^3 M.$$

For a constant step size $h > 0$ satisfying $hL < 2$, deduce that, if $y = y(x)$, then

$$e_n \leq \frac{h^3 M}{12L} \frac{1 + e^{-hL}}{1 - e^{-hL}} e^{-\frac{1}{2}hL}.$$

12.8 Show that the one-step method defined by

$$y_{n+1} = y_n + h(k_n + k_{n+1}),$$

where

$$k_n = f(x_n, y_n), \quad k_{n+1} = f(x_{n+1}, y_{n+1})$$

is consistent and has truncation error

$$T_n = -h^3 f''(f' + f' f') - (f'' + 2f'' f' + f'' f'^2) + (h^4).$$

12.9 When the classical fourth-order Runge–Kutta method is applied to the differential equation $y' = \lambda y$, where λ is a constant, show that

$$y_{n+1} = (1 + h\lambda + \frac{1}{2}h^2\lambda^2 + \frac{1}{6}h^3\lambda^3 + \frac{1}{24}h^4\lambda^4) y_n.$$

Compare this with the Taylor series expansion of $y(x_{n+1}) = y(x_n + h)$ about the point $x = x_n$.

12.10 Consider the one-step method

$$y_{n+1} = y_n + hf(x_n, y_n) + hf(x_n + h, y_n + hf(x_n, y_n)),$$

where α , β , and γ are real parameters and $h > 0$. Show that the method is consistent if, and only if, $\alpha + \beta = 1$. Show also that the order of the method cannot exceed 2.

Suppose that a second-order method of the above form is applied to the initial value problem $y' = \lambda y$, $y(0) = 1$, where λ is a positive real number. Show that the sequence (y_n) is bounded if, and only if, $h < \frac{2}{\lambda}$. Show further that, for such h ,

$$y(x_n) - y_n \leq \frac{1}{6} h^3 \lambda^3, \quad n \geq 0.$$

12.11 Find the values of a and b so that the three-step method

$$y_{n+3} + (y_{n+2} - y_n) - y_{n+1} = h(f_{n+3} + f_{n+2})$$

has order of accuracy 4, and show that the resulting method is *not* zero-stable.

12.12 Consider approximating the initial value problem $y' = f(x, y)$, $y(0) = y_0$ by the linear multistep method

$$y_{n+2} + by_{n+1} + ay_n = hf(x_n, y_n)$$

on the regular mesh $x_n = nh$ where a and b are constants.

(i) For a certain (unique) choice of a and b , this method is consistent. Find these values of a and b and verify that the order of accuracy is 1.

(ii) Although the method is consistent for the choice of a and b from part (i), the numerical solution it generates will not, in general, converge to the solution of the initial value problem as $h \rightarrow 0$, because the method is not zero-stable. Show that the method is not zero-stable for these a and b , and describe quantitatively what the unstable solutions will look like for small h .

12.13 Given that α is a positive real number, consider the linear two-step method

$$y_{n+2} - y_{n+1} = \frac{h}{3} [f(x_{n+2}, y_{n+2}) + 4f(x_{n+1}, y_{n+1}) + f(x_n, y_n)],$$

on the mesh $x_n : x_n = x_0 + nh, n = 1, 2, \dots, N$ of spacing $h, h > 0$. Determine the set of all α such that the method is zero-stable. Find α such that the order of accuracy is as high as possible; is the method convergent for this value of α ?

12.14 Which of the following linear multistep methods for the solution of the initial value problem $y' = f(x, y), y(0)$ given, are zero-stable?

- (a) $y_{n+2} - y_n = hf_n$,
- (b) $y_{n+2} + y_{n+1} - 2y_n = h(f_{n+2} + f_{n+1} + f_n)$,
- (c) $y_{n+2} - y_{n+1} = -h(f_{n+2} + 4f_{n+1} + f_n)$,
- (d) $y_{n+2} - y_n = -h(3f_{n+1} - f_n)$,
- (e) $y_{n+2} - y_n = -h(5f_{n+1} + 8f_n - f_{n-1})$.

For the methods under (a) and (c) explore absolute stability when applied to the differential equation $y' = \lambda y$ with $\lambda < 0$.

$$\begin{array}{ccc|ccc}
 -(3 & 3) & & - & & -(3 & 2 & 3) \\
 -(3+ & 3) & & -(3+2 & 3) & & - & \\
 \hline
 & & & - & & & - &
 \end{array}$$

deduce that $y' = R(h)y$, where

$$R(h) = \frac{1 + h + h^2}{1 - h + h^2}.$$

By writing $R(z)$ in the factorised form $(z+p)(z+q)/(z-p)(z-q)$, deduce that this Runge-Kutta method is A-stable.

Boundary value problems for ODEs

13.1 Introduction

In the previous chapter we discussed numerical methods for initial value problems in which all the associated side conditions for a system of differential equations are prescribed at the same point. Now we go on to consider problems where these conditions specify values at more than one point. Typically we require the solution on an interval $[a, b]$, and some conditions are given at a , and the rest at b , although more complicated situations are possible, involving three or more points.

We shall begin with the simplest case, of a second-order equation with one condition given at a and one at b . This problem is sufficient to introduce the basic ideas, and is of a type which arises quite often in practice.

We then go on to discuss the shooting method for the solution of more general problems.

13.2 A model problem

The simplest two-point boundary problem involves the second-order differential equation

$$y'' + r(x)y = f(x), \quad a < x < b, \quad (13.1)$$

with the boundary conditions

$$y(a) = A, \quad y(b) = B, \quad (13.2)$$

where A and B are given real numbers. We shall assume that r and f are given real-valued functions, defined and continuous on the bounded closed interval $[a, b]$ of the real line, and that

$$r(x) > 0, \quad a < x < b.$$

The reason for this condition will appear later, in Theorem 13.4.

We shall construct a numerical approximation to the solution on a uniform mesh of points

$$x_j = a + jh, \quad j = 0, 1, \dots, n, \quad h = (b - a)/n, \quad n \geq 2,$$

so that $x_0 = a, x_n = b$. The second derivative is approximated using the second central difference defined below.

Definition 13.1 *The central difference $\delta^2 y$ of y is defined by*

$$\delta^2 y(x) = y(x + h) - 2y(x) + y(x - h).$$

Higher-order differences are defined recursively by

$$\delta^3 y(x) = [\delta^2 y(x)] = \delta^2 y(x + h) - \delta^2 y(x - h).$$

In particular, the second central difference may be written

$$\begin{aligned} \delta^2 y(x) &= y(x + h) - 2y(x) + y(x - h) \\ &= y(x + h) - 2y(x) + y(x - h). \end{aligned}$$

Theorem 13.1 (i) *Suppose that $y \in C^4[x - h, x + h]$, i.e., that y has continuous fourth derivative on the interval $[x - h, x + h]$. Then, there exists a number ξ in $(x - h, x + h)$ such that*

$$\frac{\delta^2 y(x)}{h^2} = y''(\xi) + \frac{h^2}{12} y^{(4)}(\xi).$$

(ii) *Suppose that $y \in C^4[x - h, x + h]$; then, there exists a number ξ in $(x - h, x + h)$ such that*

$$\frac{\delta^2 y(x)}{h^2} = y''(\xi) + \frac{h^2}{12} y^{(4)}(\xi). \tag{13.3}$$

Proof (i) Taylor's Theorem shows that there exist numbers ξ_1 and ξ_2 in the intervals $(x - h, x)$ and $(x, x + h)$, respectively, such that

$$\begin{aligned} y(x - h) &= y(x) - hy'(x) + \frac{h^2}{2} y''(\xi_1) - \frac{h^3}{6} y'''(\xi_1) + \frac{h^4}{24} y^{(4)}(\xi_1), \\ y(x + h) &= y(x) + hy'(x) + \frac{h^2}{2} y''(\xi_2) + \frac{h^3}{6} y'''(\xi_2) + \frac{h^4}{24} y^{(4)}(\xi_2). \end{aligned} \tag{13.4}$$

Since y is continuous on $[x - h, x + h]$, there is a number ξ in $(x - h, x + h)$, and thus also in $(x - h, x + h)$, such that

$$-y''(\xi_1) + y''(\xi_2) = y''(\xi).$$

The required result is now obtained by adding the two equations (13.4) and dividing by h .

(ii) The proof is completely analogous, and is left to the reader as an exercise. (See Exercise 1.) \square

We can now use the central difference approximation to construct the numerical solution. Writing Y_j for the numerical approximation to $y(x_j)$, we approximate the differential equation by

$$\frac{Y_j - Y_{j-1}}{h} + r_j Y_j = f_j, \quad j = 1, 2, \dots, n-1, \quad (13.5)$$

where we have used the notation $r_j = r(x_j)$, $f_j = f(x_j)$. Now, (13.5) is a system of $n-1$ linear algebraic equations for the $n-1$ unknowns Y_j , $j = 1, 2, \dots, n-1$, with the boundary conditions specifying the values of Y_0 and Y_n ,

$$Y_0 = A, \quad Y_n = B. \quad (13.6)$$

The system may be written in matrix form as

$$M\mathbf{Y} = \mathbf{g},$$

where $\mathbf{Y} = (Y_1, \dots, Y_{n-1})^T$ and, for $n \geq 4$, the matrix M is $(n-1) \times (n-1)$ tridiagonal. Here $\mathbf{Y} = (Y_1, \dots, Y_{n-1})^T$, the nonzero elements of M are

$$M_{j,j} = 2/h + r_j, \quad M_{j,j-1} = M_{j,j+1} = -1/h, \quad (13.7)$$

and the elements of the column vector \mathbf{g} on the right-hand side are

$$g_1 = f_1 + A/h, \quad g_{n-1} = f_{n-1} + B/h, \quad g_j = f_j, \quad j = 2, 3, \dots, n-2.$$

Note how the known boundary values Y_0 and Y_n have been transferred to the right-hand side, and appear in the first and last elements of \mathbf{g} . The solution of this system is very easy, using the algorithm for tridiagonal matrices described in Section 3.3. Using the fact that $r(x) \geq 0$, we see that the off-diagonal elements of M are negative, the diagonal elements are positive, and in each row the diagonal element is at least as large as the sum of absolute values of the off-diagonal elements. Theorem 3.4 implies that no row interchanges are needed in the calculation, and that the matrix M is nonsingular. The calculation is therefore very straightforward and efficient, and requires very little computational time, even for a mesh which may contain several hundred points.

Now, $y(x_j)$ and Y_j satisfy

$$\begin{aligned} L(y(x_j)) &= f + T_j, \quad j = 1, 2, \dots, n-1, \\ L(Y_j) &= f, \quad j = 1, 2, \dots, n-1, \end{aligned}$$

from the definition of truncation error and (13.5); hence, by subtraction,

$$L(e_j) = T_j, \quad j = 1, 2, \dots, n-1,$$

with the boundary conditions $e_0 = e_n = 0$. We must now use the bound on T_j to derive a bound on the error e_j . This will be achieved by means of the following theorem.

Theorem 13.3 (Maximum Principle) *Let $a_j, b_j, c_j, j = 0, 1, \dots, n$, be positive real numbers such that $b_j = a_j + c_j$, and suppose that $u_j, j = 0, 1, \dots, n$, are real numbers such that*

$$a_j u_{j-1} + b_j u_j - c_j u_{j+1} \leq 0, \quad j = 1, 2, \dots, n-1.$$

Then, $u_j \leq K, j = 0, 1, \dots, n$, where $K = \max\{u_0, u_n, 0\}$.

Proof Let $u = \max\{u_0, u_1, \dots, u_n\}$; then if $r = 0, r = n$, or $u = 0$ the result is trivial. Suppose then that $1 \leq r \leq n-1$, and that $u_r > 0$. Since u_r is the maximum of the u_j , we know that

$$u_r \geq u_{r-1}, \quad u_r \geq u_{r+1}.$$

Hence

$$\begin{aligned} b_r u_r &\leq a_r u_{r-1} + c_r u_{r+1} \\ &\leq a_r u_r + c_r u_r \\ &= b_r u_r, \end{aligned}$$

since $u_r > 0$. This means that equality holds throughout, so that $u_{r-1} = u_r = u_{r+1}$. We can then apply the same argument to both u_{r-1} and u_{r+1} , continuing until we find that either $u = u_0$ or $u = u_n$. Thus, in this case $u = u_r = \max\{u_0, u_1, \dots, u_n\}$, as required. \square

Theorem 13.4 *Suppose that the solution y of the boundary value problem (13.1), (13.2) has a continuous fourth derivative on $[a, b]$, and that $Y_j, j = 0, 1, \dots, n$, is the solution of the central difference approximation (13.5), (13.6). Then,*

$$\max |y(x_j) - Y_j| \leq h^4 (b-a) M. \tag{13.10}$$

Proof Let $e = y(x) - Y$. We have already seen that $L(e) = T$, $j = 1, 2, \dots, n - 1$. Defining

$$e = C \sum_{j=0}^{n-1} (2j - n) h^{2j}, \quad j = 0, 1, \dots, n, \quad (13.11)$$

where C is a constant, we see that

$$\begin{aligned} L(e) &= C \sum_{j=1}^{n-1} (2j - 2 - n) \cdot 2(2j - n) + (2j + 2 - n)^2 + r \\ &= 8C + r, \quad j = 1, 2, \dots, n - 1. \end{aligned}$$

Hence

$$L(e + Y) = T + 8C + r, \quad j = 1, 2, \dots, n - 1.$$

If we choose $C = T/8$ with $T = -h M$, we see that $L(e + Y) \leq 0$, since $T \leq T$, $r \leq 0$ and 0 , and L satisfies the conditions of the Maximum Principle. Now,

$$e + Y = e + Y = 0,$$

so that, according to Theorem 13.3, $e + Y \leq 0$ for $j = 0, 1, \dots, n$. However, $Cn h^{2n} \leq 0$, so we have the result

$$e \leq Cn h^{2n} = -(b - a) T = -h (b - a) M, \quad j = 0, 1, \dots, n.$$

By applying the same argument to $L(e - Y)$ we find that

$$e \leq -h (b - a) M, \quad j = 0, 1, \dots, n.$$

Combining these upper bounds for e and $e - Y$ gives the required result. □

The function e defined by (13.11) is called a **comparison function**. An alternative proof of Theorem 13.4, based on the properties of monotone matrices, can be given by using the result in Exercise 2. Notice that the condition $r(x) \leq 0$ is used in the application of the Maximum Principle in the above proof.

This theorem shows that, provided the solution y has a continuous fourth derivative, the numerical method is **convergent**, that is

$$\max y(x) - Y \leq 0 \quad \text{as } n \rightarrow \infty$$

(or, equivalently, as $h = (b - a)/n \rightarrow 0$). This means that we can obtain any required accuracy by choosing n sufficiently large.

Note that the approximation to $y(x)$ at $x = x_j$ may be written

$$\frac{-[y(x_{j+1/2}) + y(x_{j-1/2})]}{h}.$$

For $j = 1, 2, \dots, n - 1$, we define the truncation error T_j as in Definition 13.2. In addition, since we shall now also incur an error in the approximation of the boundary condition at $x = a$, we define

$$T_0 = \frac{2(1 + \alpha h)}{h} + r y(0) - \frac{2}{h}y(h) - f + \frac{2}{h}A.$$

The aim of our next result is to quantify the size of the truncation error in terms of the mesh size h .

Theorem 13.6 *Suppose that the solution y to the boundary value problem (13.1), (13.2) has a continuous fourth derivative on the closed interval $[a - h, b]$. Then, the truncation error of the central difference approximation to (13.1) with boundary conditions (13.12) may be written*

$$\begin{aligned} T_j &= -\frac{h^4}{12} y^{(4)}(\xi_j), \quad j = 1, 2, \dots, n - 1, \\ T_0 &= -\frac{h^4}{12} y^{(4)}(\xi_0) - \frac{h^4}{24} y^{(4)}(\xi_1), \end{aligned}$$

for some value of ξ_j in the interval $(x_{j-1/2}, x_{j+1/2})$, $1 \leq j \leq n - 1$, and some value ξ_0 in the interval $(x_{-1/2}, x_{1/2})$ where $x_{-1/2} = a - h$.

Proof For $j = 1, 2, \dots, n - 1$, this is the same result as in Theorem 13.2. When $j = 0$, we find that

$$\begin{aligned} T_0 &= \frac{2(1 + \alpha h)}{h} + r y(0) - \frac{2}{h}y(h) - f + \frac{2}{h}A \\ &= \frac{y(h) - 2y(0) + y(-h)}{h} + r(0)y(0) - f(0) \\ &\quad + \frac{2}{h} \frac{y(h) - y(-h)}{2h} - y(0) - A \\ &= -\frac{h^4}{12} y^{(4)}(\xi_0) - \frac{h^4}{24} y^{(4)}(\xi_1), \end{aligned}$$

where we have used Theorem 13.5. □

Theorem 13.7 *Suppose that the solution y of (13.1) with the boundary conditions (13.12) has a continuous fourth derivative on the interval $[a - h, b]$; then, the numerical solution obtained from the central difference*

approximation satisfies

$$\max_{y(x)} |y - Y| \leq \frac{h}{2} \left[(b-a)M + (b-a)M^2 \right].$$

Proof The proof is very similar to that of Theorem 13.4, but requires the use of a more complicated comparison function. Let us define

$$L(u) = \frac{u}{h} + r u$$

In particular,

$$e + \max e + ,0 . \tag{13.15}$$

However, $L(e +) = 0$; thus, by the definition of $L(e +)$,

$$e + \frac{2}{2(1 + h) + hr} (e +) .$$

On writing $= 2/(2(1 + h) + hr)$ and noting that, since > 0 and $r = 0$, we have $0 < < 1$, it follows that

$$e + (e +) . \tag{13.16}$$

Inserting this inequality into the left-hand side of (13.15), we find that

$$e + \max (e +),0 .$$

If $e +$ were positive, this inequality and the fact that $0 < < 1$ would imply $e + = 0$, leading to a contradiction. Therefore, $e + = 0$. Returning with this information to (13.14), we conclude that $e + = 0$ for $j = 0, 1, \dots, n$, and the rest of the proof then follows as in the proof of Theorem 13.1. □

13.5 The general self-adjoint problem

The general self-adjoint boundary value problem is

$$\frac{d}{dx} p(x) \frac{dy}{dx} + r(x)y = f(x), \quad a < x < b, \tag{13.17}$$

where r and f are real-valued functions, defined and continuous on $[a, b]$, p is a real-valued continuously differentiable function on $[a, b]$, $r(x) \geq 0$ and $p(x) = c > 0$. We shall consider only the case where the boundary conditions prescribe the values of y at each end,

$$y(a) = A, \quad y(b) = B . \tag{13.18}$$

The central difference approximation to the equation (13.17) may be written

$$\frac{(p_j Y_j)}{h} + r_j Y_j = f_j, \quad j = 1, 2, \dots, n - 1,$$

or, in detail,

$$\frac{p_j (Y_j - Y_{j-1}) + p_{j+1} (Y_{j+1} - Y_j)}{h} + r_j Y_j = f_j, \tag{13.19}$$

for $j = 1, 2, \dots, n$, and is supplemented by the boundary conditions

$$Y_0 = A, \quad Y_n = B. \tag{13.20}$$

It is easy to see that this represents a system of linear equations for the unknowns Y_1, Y_2, \dots, Y_{n-1} , and that the matrix of the system is tridiagonal and diagonally dominant, just as it was in the special case (13.1), which corresponds to $p(x) = 1$. The solution of the system is therefore a very simple matter.

Next, we consider the error analysis of the difference scheme (13.19), (13.20). We begin by quantifying the size of the truncation error

$$T_j = \frac{(p_j y(x_j))}{h} + r_j y(x_j) - f_j, \quad j = 1, 2, \dots, n-1,$$

in terms of the mesh size h .

Lemma 13.1 *Suppose that $p \in C^1[a, b]$ and $y \in C^2[a, b]$. The truncation error T_j of the central difference approximation (13.19) then satisfies*

$$|T_j| \leq h^2 \max_{x \in [a, b]} \left(\frac{1}{2} |p''(x)| + |p'(x)| + 2|p(x)| \right) |y''(x)|,$$

for $j = 1, 2, \dots, n-1$.

Proof By expanding in Taylor series as we have done before, we find that

$$\begin{aligned} p_j [y(x_{j+1}) - y(x_j)] &= p_j [hy_{j+1/2} + \frac{1}{2}h^2 y_{j+3/2} + \dots], \\ p_{j-1} [y(x_j) - y(x_{j-1})] &= p_{j-1} [hy_{j-1/2} - \frac{1}{2}h^2 y_{j-3/2} + \dots], \end{aligned}$$

where $(x_{j+1/2}, x_{j+3/2})$ and $(x_{j-1/2}, x_{j-3/2})$. The first term in the difference of these expressions gives, in the same way,

$$h[p_j y(x_{j+1/2}) - p_{j-1} y(x_{j-1/2})] = h[h(p_j y_{j+1/2}) - h(p_{j-1} y_{j-1/2})]$$

where $(x_{j-1/2}, x_{j+1/2})$. For the other term we can write

$$\begin{aligned} -h p_j y_{j+1/2} + p_{j-1} y_{j-1/2} &= -h \left(\frac{1}{2} p_j y_{j+3/2} - \frac{1}{2} p_{j-1} y_{j-3/2} \right) + p_{j-1} [y_{j+1/2} - y_{j-1/2}] \\ &\quad - h \left(\frac{1}{2} p_j y_{j+1/2} + p_{j-1} y_{j-1/2} \right), \end{aligned}$$

since $h < 2h$. Here, $(x_{j-1/2}, x_{j+1/2})$ and $(x_{j+1/2}, x_{j-1/2})$ lies between $x_{j-1/2}$ and $x_{j+1/2}$. The required bound follows immediately. \square

that $Y_j, j = 0, 1, \dots, n$, is the solution of the central difference approximation (13.19), (13.20). Then, with T as in Lemma 13.1,

$$\max |y(x) - Y_j| \leq \frac{1}{4} T. \tag{13.21}$$

Proof The proof of this theorem follows that of Theorem 13.4, using the bound from Lemma 13.1 on the truncation error and the comparison function from Lemma 13.2. The details are left as an exercise. \square

13.6 The Sturm–Liouville eigenvalue problem

Suppose that r is a real-valued function, defined and continuous on the closed interval $[a, b]$, p is a real-valued function, defined and continuously differentiable on $[a, b]$, and $r(x) \geq 0, p(x) > 0$ for all $x \in [a, b]$. The differential equation

$$\frac{d}{dx} \left(p(x) \frac{dy}{dx} \right) + r(x)y = \lambda y, \quad a < x < b, \tag{13.22}$$

with homogeneous boundary conditions $y(a) = y(b) = 0$, has only the trivial solution $y = 0$, except for an infinite sequence of positive *eigenvalues* $\lambda_m = \lambda_m$, $m = 1, 2, \dots$. We shall now consider a numerical method for finding these eigenvalues and the corresponding *eigenfunctions*, $y_m(x)$, $m = 1, 2, \dots$

In the simple case where $p(x) = 1$ and $r(x) = 0$ the solution to this problem is, of course, $\lambda_m = [m\pi/(b-a)]^2, y_m(x) = A \sin m\pi t$, $m = 1, 2, \dots$, where A is a nonzero constant and $t = (x-a)/(b-a)$.

Using the same finite difference approximation as in the previous section, we obtain the equations

$$\frac{p_j (Y_{j+1} - Y_j) - p_{j-1} (Y_j - Y_{j-1})}{h} + r_j Y_j = \lambda Y_j, \quad j = 1, 2, \dots, n-1.$$

Together with the boundary conditions $Y_0 = Y_n = 0$, this shows that λ is an eigenvalue of a symmetric tridiagonal matrix M whose entries are

$$M_{jj} = \frac{p_j + p_{j-1}}{h} + r_j, \quad 1 \leq j \leq n-1, \\ M_{j,j-1} = M_{j-1,j} = -\frac{p_{j-1}}{h}, \quad 2 \leq j \leq n, \quad M_{11} = M_{nn} = \frac{p_1}{h}, \quad 1 \leq j \leq n-1,$$

and the approximate function values Y_j are the elements of the corresponding eigenvector. This algebraic eigenvalue problem is easily solved by the method described in Chapter 5.

The boundary value problems which we have discussed so far have all had a unique solution. The eigenvalue problem (13.22) has an infinite number of solutions, and the mesh used in the numerical computation has to be chosen to adequately represent the eigenfunctions required – the computation can obviously only find a finite number of them. The matrix M has $n - 1$ eigenvalues and eigenvectors and, as we shall see, it will normally give a good approximation to the first few eigenvalues, $\lambda_1, \lambda_2, \dots$, and a much less accurate approximation to λ_{n-1} .

To analyse the error in the eigenvalue we proceed as before, by defining the truncation error

$$T_j = \frac{p(x_j)(y_j - y_{j-1}) - p(x_{j-1})(y_{j-1} - y_{j-2})}{h} + r(x_j)y_j - y_{j-1},$$

$j = 1, 2, \dots, n - 1,$

where $y_j = y(x_j)$. These equations can now be written

$$\begin{aligned} (M - \Lambda I)\mathbf{Y} &= \mathbf{0}, \\ (M - I)\mathbf{Y} &= \mathbf{T}, \end{aligned}$$

where

$$\begin{aligned} \mathbf{Y} &= (Y_1, \dots, Y_{n-1})^T, \\ &= (y_1, \dots, y_{n-1})^T, \\ \mathbf{T} &= (T_1, \dots, T_{n-1})^T. \end{aligned}$$

Theorem 5.15 of Chapter 5 applies to this problem, and shows that one of the eigenvalues, Λ , of the matrix M satisfies

$$\Lambda = \mathbf{T}^T / \mathbf{Y}^T. \quad (13.23)$$

In the simpler case where $p(x) = 1$ and $r(x) = 0$ the truncation error is

$$T_j = -h y''(x_j), \quad (x_{j-1}, x_{j+1}),$$

so the numerical method has evaluated the eigenvalue with error less than

$$-h \max_{j=1, \dots, n-1} |y''(x_j)| = -h \max_{j=1, \dots, n-1} |y''(x_j)|.$$

Since the m th eigenfunction $y_m(x)$ is given by

$$y_m(x) = y_m(x) = \sin(m(x-a)/(b-a)), \quad x \in (a, b),$$

we see that

$$y_m(x) = \frac{m}{b-a} y(x), \quad x \in (a, b).$$

This shows that, for example, the error in the tenth eigenvalue, corresponding to $m = 10$, is likely to be about 10 times larger than the error in the first eigenvalue; more generally, to evaluate higher eigenvalues of the equation will require the use of a smaller interval h .

13.7 The shooting method

The methods we have described for the linear boundary value problem may be extended to nonlinear differential equations. We shall not discuss how this is done; instead, we shall describe an alternative approach, called the **shooting method**. We shall consider the nonlinear model problem

$$y' = f(x, y), \quad a < x < b, \quad y(a) = A, \quad y(b) = B,$$

where we assume that the function $f(x, y)$ is continuous and differentiable, and that

$$\frac{\partial f}{\partial y}(x, y) > 0, \quad a < x < b, \quad y \in \mathbb{R}.$$

The central idea of the method is to replace the boundary value problem under consideration by an initial value problem of the form

$$y' = f(x, y), \quad a < x < b, \quad y(a) = A, \quad y(b) = t,$$

where t is to be chosen in such a way that $y(b) = B$. This can be thought of as a problem of trying to determine the angle of inclination $\tan^{-1} t$ of a loaded gun, so that, when shot from height A at the point $x = a$, the bullet hits the target placed at height B at the point $x = b$. Hence the name, shooting method.

Once the boundary value problem has been transformed into such an ‘equivalent’ initial value problem, any of the methods for the numerical solution of initial value problems discussed in Chapter 12 can be applied to find a numerical solution. Thus, in particular, the costly exercise of solving a large system of nonlinear equations, arising from a direct finite

difference approximation of the nonlinear boundary value problem, can be completely avoided.

If we write

$$y(a) = t,$$

a numerical solution of the differential equation with the initial conditions $y(a) = A, y(a) = t$ can be obtained by any of the methods of Chapter 12. This solution will depend on t , and we may write it as $y(x; t)$. In particular the value at $x = b$ will be a function of t ,

$$y(b; t) = \phi(t). \tag{13.24}$$

The solution of the nonlinear boundary value problem therefore reduces to the determination of the value of t for which the boundary condition at $x = b$ is also satisfied, *i.e.*,

$$\phi(t) - B = 0.$$

There are a number of well-known methods for the solution of equations of this form; Newton's method is an obvious example. Generally, we shall not, of course, have a closed form expression for the function $\phi(t)$, in general, but this is not necessary; all that is needed is a numerical algorithm to calculate the value of $\phi(t)$ for a given value of t , and this we have. To use Newton's method we shall also need to be able to calculate the value of $\phi'(t)$, and this is easily done.

The function $y(x; t)$ is defined, for all t , as the solution of the initial value problem

$$y'(x; t) = f(x, y(x; t)), \quad y(a; t) = A, \quad y(a; t) = t, \tag{13.25}$$

where $\frac{d}{dx}$ and $\frac{d}{dt}$ indicate differentiation with respect to the variable x . We can differentiate these throughout with respect to t , giving

$$-\frac{dy}{dt}(x; t) = \frac{df}{dy}(x, y(x; t)) \frac{y}{dt}(x; t), \quad \frac{y}{dt}(a; t) = 0, \quad \frac{y}{dt}(a; t) = 1.$$

Writing

$$w(x, t) = -\frac{y}{dt}(x; t),$$

and interchanging the order of differentiation, we find that $w(x; t)$ may be obtained as the solution of the initial value problem

$$w'(x; t) = w(x; t) \frac{df}{dy}(x, y(x; t)), \quad w(a; t) = 0, \quad w(a; t) = 1. \tag{13.26}$$

By virtue of (13.24), the required derivative is then given by

$$\dot{\mathbf{t}} = \mathbf{w}(\mathbf{b}, \mathbf{t}).$$

To implement this method, it is convenient to solve the two initial value problems, (13.25) and (13.26), in tandem, by writing them as a system of four simultaneous first-order differential equations:

Proof Suppose that the solution of the system of differential equations with $t = t_j$ is $u(x; t_j)$, $i = 1, 2, 3, 4$, and the corresponding numerical solution is $v_j(t_j)$, $i = 1, 2, 3, 4$, $j = 1, 2, \dots, n$; then

$$u(x; t_j) - v_j(t_j) = C(t_j)h^j.$$

Moreover $v_j(t_j) \in B_j$, so that

$$u(b; t_j) \in B_j = u(b; t_j) - v_j(t_j) + v_j(t_j) \in B_j + C(t_j)h^j. \tag{13.28}$$

Let us write $u(x; t) = y(x) + u(x; t)$; by subtraction we see that

$$\begin{aligned} u(x; t) &= y(x) + u(x; t) \\ &= f(x, y(x)) - f(x, u(x; t)) \\ &= (x; t) \frac{f}{y}(x, u(x; t)), \end{aligned}$$

where $u(x; t)$ lies between $u(x; t)$ and $y(x)$.

Suppose that $u(a; t) > 0$; since $u(a; t) = 0$, there is some interval to the right of a in which $u(x; t) > 0$. Then, either $u(x; t) > 0$ for the whole of $[a, b]$, or there is a value c such that $a < c < b$ and $u(c; t) = 0$. In the latter case, $u(x; t)$ must vanish at some point $x = d$ between a and c . However, in the interval $[a, d]$, $u(x; t) > 0$ and $f/y > 0$, so that $u(x; t) > 0$. Consequently, in the interval $[a, d]$, $u(x; t) > u(a; t) > 0$, and we have a contradiction. Thus, $u(x; t) > 0$ for all $a < x \leq b$. It then follows that $u(x; t)$, and hence also $y(x) + u(x; t)$ are positive on the whole interval $[a, b]$, which means that $x \mapsto y(x) + u(x; t)$ is monotonic increasing on $[a, b]$. If we had begun with the assumption that $u(a; t) < 0$ an analogous argument shows that $x \mapsto y(x) + u(x; t)$ would have been monotonic decreasing on $[a, b]$. It is left to the reader to discuss the trivial case when $u(a; t) = 0$.

In any case,

$$u(x; t) \leq u(b; t), \quad a \leq x \leq b,$$

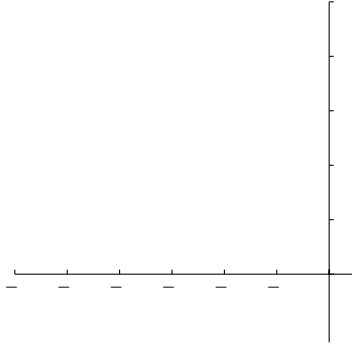
and therefore, since $y(b) = B$ and recalling (13.28),

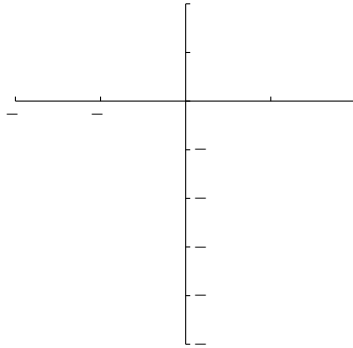
$$y(x) + u(x; t) \leq B + u(b; t) = C(t)h^j + B.$$

Thus, finally,

$$y(x) - v_j(t_j) = y(x) + u(x; t_j) - u(x; t_j) - v_j(t_j) \leq C(t_j)h^j + B - v_j(t_j) = C(t_j)h^j + C(t_j)h^j, \quad j = 1, 2, \dots, n,$$

and hence the desired bound. □





& 5 , B , * & \$+

13.8 Notes

The following books are standard texts on the subject of numerical approximation of boundary value problems:

! , , *Numerical Methods for Two-Point Boundary Value Problems*, Reprint of the 1968 original published by Blaisdell, Dover, New York, 1992.

! , , *Numerical Solution of Two-Point Boundary Value Problems*, SIAM, Philadelphia, fourth printing, 1990.

A more recent survey of the subject is found in

O +) " % + + + 4 % % " *Numerical Solution of Boundary Value Problems for Ordinary Differential Equations*, Corrected reprint of the 1988 original, Classics in Applied Mathematics, 13, SIAM, Philadelphia, 1995.

In practical implementations of the shooting method into mathematical software (see, for example, Appendix A in the Ascher *et al.* book), the interval $[a, b]$ is subdivided into smaller intervals on each of which the shooting method is applied with appropriately chosen initial values. The ‘initial’ conditions on the subintervals are then simultaneously adjusted in order to satisfy the boundary conditions and appropriate continuity conditions at the points of the subdivision. From the practical viewpoint, this extension of the basic shooting method considered in this chapter is extremely important: the various difficulties which may arise in the implementation of the basic method (such as, for example, growth of the

solution to the initial value problem over the interval $[a, b]$, leading to loss of accuracy in the solution of the equation $y'(t) = B y(t)$ are discussed, for example, in Section 2.4 of the 1992 book by Keller.

Sturm–Liouville problems originated in a paper of Jacques Charles François Sturm: Sur les équations différentielles linéaires du second ordre, *J. Math. Pures Appl.* **1**, 106–186, 1836, in Joseph Liouville’s newly founded journal. Sturm’s paper was followed by a series of articles by Sturm and Liouville in subsequent volumes of the journal. They examined general linear second-order differential equations, the properties of their eigenvalues, the behaviour of the eigenfunctions and the series expansion of arbitrary functions in terms of these eigenfunctions. An extensive survey of the theory and numerical analysis of Sturm–Liouville problems can be found in

§ " *Numerical Solution for Sturm–Liouville Problems*,
Oxford University Press Monographs in Numerical Analysis, Clarendon Press, Oxford, 1993.

See also Section 11.3, page 478, of the Ascher *et al.* book cited above.

Exercises

- 13.1 Suppose that $y \in C[x-h, x+h]$; show that there exists a real number $\theta \in (x-h, x+h)$ such that

$$\frac{y(x)}{h} = y(x) + -h y'(x) + \frac{1}{2} h^2 y''(\theta).$$

- 13.2 Use Theorem 3.6 to show that the matrix M in (13.7) is monotone. Use the result of Exercise 4 to show that $M^{-1} \leq -$.
- 13.3 On the interval $[a, b]$ the differential equation

$$y' + f(x)y = g(x)$$

is approximated by

$$\frac{y}{h} + - y - + y + y = - g - + g + g ,$$

where $-$, $,$ and $+$ are constants. Assuming that the solution y has the appropriate number of continuous derivatives, show that the truncation error of this approximation may be written as follows:

(i) if $\alpha + \beta + \gamma = 1$, then

$$T = (\alpha + \beta + \gamma)y(x) + Z h,$$

where $Z = (\alpha + \beta + \gamma)M$;

(ii) if $\alpha + \beta + \gamma = 1$ and $\alpha = 0$, then

$$T = (\beta - \gamma)hy(x) + Z h,$$

where $Z = [-(\beta + \gamma) + \alpha]M$;

(iii) if $\alpha + \beta + \gamma = 1$, $\alpha = 0$ and $\gamma = -1$, then

$$T = (\beta - \gamma)hy(x) + Z h,$$

where $Z = [-\beta + \gamma]M$;

(iv) if $\alpha = 0$ and $\gamma = -1$, then

$$T = -hy(x) + Z h,$$

where $Z = -M$.

13.4 The approximation of Exercise 3 is used, with the values $\alpha = 0$, $\beta = 1/12$, $\gamma = 5/6$. Use Taylor's Theorem with integral remainder (Appendix, Theorem A.5) to show that the truncation error of this approximation may be written

$$T = \int_{-\alpha}^{\beta} G(s)y(x+s)ds,$$

where

$$G(s) = (h-s)^5/5! - h(h-s)^3/3!, \quad 0 \leq s \leq h,$$

with a similar expression for $h-s \leq 0$. Show that $G(s) \leq 0$ for all $s \in [-h, h]$, and hence use the Integral Mean Value Theorem to show that the truncation error can be expressed as

$$T = \frac{h}{240}y(\xi)$$

for some value of ξ in $(x-h, x+h)$.

- 13.5 Suppose that the solution of (13.1), (13.2) has a continuous sixth derivative on $[a, b]$, and that Y is the solution of the approximation used in Exercise 4. Show that

$$y(x) - Y \approx -\frac{h^6}{720} y^{(6)}(\xi), \quad j = 0, \dots, n,$$

provided that

$$h r(x) \leq 12, \quad j = 1, \dots, n-1.$$

- 13.6 Complete the proof of Theorem 13.7.
 13.7 Show that the solution of the boundary value problem

$$y'' + a^2 y = 0, \quad y(-1) = 1, \quad y(1) = 1,$$

is

$$y(x) = \frac{\cosh ax}{\cosh a}.$$

Use the identity

$$\cosh(x+h) + \cosh(x-h) = 2 \cosh x \cosh h$$

to verify that the solution of the difference approximation (13.5) to this problem is

$$Y = \frac{\cosh x}{\cosh a},$$

where

$$= (1/h) \cosh^{-1}(1 + ah).$$

By expanding in Taylor series, show that

$$Y = y(x) + \frac{h^2 a^2}{24} (\cosh ax \sinh a - x \sinh ax \cosh a) / (\cosh a) + O(h^4).$$

Verify that this result is consistent with Theorem 13.4 when h is small.

- 13.8 Carry out a similar analysis as in Exercise 7 for the boundary value problem

$$y'' - a^2 y = 0, \quad y(0) = 0, \quad y(1) = 1,$$

and explain why in this case Theorem 13.4 cannot be used. What restriction is required on the value of a ?

13.9 The eigenvalue problem

$$y'' = \mu y, \quad y(0) = y(1) = 0,$$

is approximated by

$$\frac{Y_j - 2Y_{j-1} + Y_{j-2}}{h^2} = \mu Y_{j-1}, \quad j = 1, \dots, n-1, \quad Y_0 = Y_n = 0.$$

Show that the differential equation has solution $y = \sin m\pi x$, $\mu = m^2\pi^2$ for any positive integer m . Show also that the difference approximation has solution $Y_j = \sin m\pi x_j$, $j = 0, 1, \dots, n$, and give an expression for the corresponding value of μ . Use the fact that

$$1 - \cos \theta = 2 \sin^2 \frac{\theta}{2},$$

to show that $\mu = m^2\pi^2 h^2/12$, and compare with the bound given by (13.23).

The finite element method

14.1 Introduction: the model problem

In Chapter 13 we explored finite difference methods for the numerical solution of two-point boundary value problems. The present chapter is devoted to the foundations of the theory of finite element methods. For the sake of simplicity the exposition will be, at least initially, confined to the second-order ordinary differential equation

$$\frac{d}{dx} p(x) \frac{du}{dx} + r(x)u = f(x), \quad a < x < b, \quad (14.1)$$

where $p \in C[a, b]$, $r \in C[a, b]$, $f \in L^2(a, b)$ and $p(x) > 0$, $r(x) \geq 0$ for all $x \in [a, b]$, subject to the boundary conditions

$$u(a) = A, \quad u(b) = B. \quad (14.2)$$

Later on in the chapter, in Section 14.5, we shall also consider the ordinary differential equation

$$\frac{d}{dx} p(x) \frac{du}{dx} + q(x) \frac{du}{dx} + r(x)u = f(x), \quad a < x < b, \quad (14.3)$$

subject to the boundary conditions (14.2). Indeed, much of the material discussed here can be extended to partial differential equations; for pointers to the relevant literature we refer to the Notes at the end of the chapter.

The finite element method was proposed in a paper by Richard Courant in the early 1940s, although the historical roots of the method can be traced back to earlier work by Galerkin in 1915; unfortunately, the relevance of Courant's article was not recognised at the time and the idea was forgotten. In the early 1950s the method was rediscovered by engineers, but its systematic mathematical analysis began only a decade later. Since then, the finite element method has been developed into one of the most general and powerful techniques for the numerical solution of differential equations which is widely used in engineering design and analysis.

Unlike finite difference schemes which seek to approximate the unknown analytical solution to a differential equation at a finite number of selected points, the grid points or mesh points in the computational domain, the finite element method supplies an approximation to the analytical solution in the form of a piecewise polynomial function, defined over the entire computational domain. For example, in the case of the boundary value problem (14.1), (14.2), the simplest finite element method uses a linear spline, defined over the interval $[a, b]$, to approximate the analytical solution u .

We shall consider two techniques for the construction of finite element approximations: the **Rayleigh–Ritz principle** and the **Galerkin principle**. In the case of the boundary value problem (14.1), (14.2) the approximations which stem from these two principles will be seen to coincide. We note, however, that since the Rayleigh–Ritz principle relies on the fact that the boundary value problem under consideration can be restated as a variational problem involving the minimisation of a certain quadratic functional over a function space, its use is restricted to *symmetric* boundary value problems, such as (14.1), (14.2) where (14.1) does not contain a first-derivative term; for example, the Rayleigh–Ritz principle is not applicable to (14.3), (14.2) unless $q(x) = 0$. The precise sense in which the word *symmetric* is to be interpreted here will be clar-

1 4 (1 1 G B. K 6
 5 ' #C * # @ * L , D C & K 6 # @ &
 # . . . : + ! " % D 4 9 " 1 ' , 2 ; <
 / 5 O " ! " # @ . ? 4
 2 F 5, F " B ' , # . & # " B D C #
 K # @ > ' , D 4 / " 1 6 5 \$
 , ,) 4 O 8 1 6 5
 1 G # @ # > 4 9 # @ 8 F 1 \$
 1) 1 ' , , 1 5 % , 6 1 , , 4

14.2 Rayleigh–Ritz and Galerkin principles

The Rayleigh–Ritz principle relies on converting the boundary value problem (14.1), (14.2) into a variational problem involving the minimisation of a certain quadratic functional over a function space.

Let us define the quadratic functional $J : H_0^1(a, b) \rightarrow \mathbb{R}$ by

$$J(w) = \frac{1}{2} \int_a^b [p(x)(w')^2 + r(x)w^2] dx - \int_a^b f(x)w(x) dx$$

where $w \in H_0^1(a, b)$, and consider the following *variational problem*:

$$(RR) \quad \text{find } u \in H_0^1(a, b) \text{ such that } J(u) = \min_{w \in H_0^1(a, b)} J(w),$$

which we shall henceforth refer to as the **Rayleigh–Ritz principle**. For the sake of notational simplicity we define

$$(w, v) = \int_a^b [p(x)w'(x)v'(x) + r(x)w(x)v(x)] dx$$

and recall from Chapter 9 the definition of inner product on $L^2(a, b)$:

$$(w, v) = \int_a^b w(x)v(x) dx.$$

Using these, we can rewrite $J(w)$ as follows:

$$J(w) = \frac{1}{2} (w, w) - (f, w), \quad w \in H_0^1(a, b). \tag{14.4}$$

The mapping $J : H_0^1(a, b) \rightarrow \mathbb{R}$ is a **bilinear functional** in the following sense:

$$\begin{aligned} J(\alpha w + \beta v) &= \alpha J(w) + \beta J(v) \\ &\text{for all } \alpha, \beta \in \mathbb{R} \text{ and all } w, v \in H_0^1(a, b); \\ J(w, \mu v + \nu w) &= \mu J(w, v) + \nu J(w, w) \\ &\text{for all } \mu, \nu \in \mathbb{R} \text{ and all } w, v \in H_0^1(a, b). \end{aligned}$$

We note, in addition, that the bilinear functional $J(\cdot, \cdot)$ is **symmetric**, in that

$$J(w, v) = J(v, w) \quad w, v \in H_0^1(a, b). \tag{14.5}$$

Our next result provides an equivalent characterisation of the Rayleigh–Ritz principle; it relies on the fact that the bilinear functional $J(\cdot, \cdot)$ is symmetric in the sense of (14.5).

Theorem 14.1 A function u in $H_{\mathfrak{g}}(a, b)$ minimises $+(\)$ over $H_{\mathfrak{g}}(a, b)$ if, and only if,

$$(G) \quad (u, v) = f, v! \quad v \in H(a, b). \quad (14.6)$$

This identity will be referred to as the **Galerkin principle**.

Proof of theorem Suppose that $u \in H_{\mathfrak{g}}(a, b)$ minimises $+(\)$ over $H_{\mathfrak{g}}(a, b)$; that is, $+(u) \leq +(w)$ for all $w \in H_{\mathfrak{g}}(a, b)$. Noting that $w = u + v$ belongs to $H_{\mathfrak{g}}(a, b)$ for all u and all $v \in H(a, b)$, we deduce that

$$\begin{aligned} +(u) &\leq +(u + v) = -((u + v, u + v) - f, u + v! \\ &= +(u) + [((u, v) - f, v!] + - ((v, v) \end{aligned} \quad (14.7)$$

for all $v \in H(a, b)$ and all u . Here, in the transition from the first line to the second we made use of the fact that $((u, v) = ((v, u)$ for all v in $H(a, b)$, which follows from (14.5). Now, (14.7) implies that

$$- ((v, v) \leq [((u, v) - f, v!]$$

for all $v \in H(a, b)$ and all u . Let us suppose that $\epsilon > 0$, divide both sides of the last inequality by ϵ and pass to the limit $\epsilon \rightarrow 0$ to deduce that

$$0 \leq ((u, v) - f, v! \quad v \in H(a, b). \quad (14.8)$$

On replacing v by $-v$ in (14.8), we have that also

$$0 \leq ((u, v) - f, v! \quad v \in H(a, b). \quad (14.9)$$

We conclude from (14.8) and (14.9) that

$$((u, v) = f, v! \quad v \in H(a, b), \quad (14.10)$$

as required.

Conversely, if $u \in H_{\mathfrak{g}}(a, b)$ is such that $((u, v) = f, v!$ for all v in $H(a, b)$, then

$$+(u + v) = +(u) + [((u, v) - f, v!] + - ((v, v) + (u)$$

for all $v \in H(a, b)$ and all u ; therefore, u minimises $+(\)$ over $H_{\mathfrak{g}}(a, b)$. \square

Thus we have shown that, as long as $((\ , \)$ is a symmetric bilinear functional, $u \in H_{\mathfrak{g}}(a, b)$ satisfies the Rayleigh–Ritz principle if, and only if, it satisfies the Galerkin principle. Our next task is to explain the

relationship between (RR) and (G) on the one-hand and (14.1), (14.2) on the other. Since in the case of a symmetric bilinear functional (\cdot, \cdot) the principles (RR) and (G) are equivalent, it is sufficient to clarify the connection between (G), for example, and the boundary value problem (14.1), (14.2).

We begin with the following definition.

Definition 14.3 *If a function $u \in H_0(a, b)$ satisfies the Galerkin principle (14.6), it is called a **weak solution** to the boundary value problem (14.1), (14.2), and the Galerkin principle is referred to as the **weak formulation** of the boundary value problem (14.1), (14.2).*

Let us justify this terminology. Suppose that $u \in H_0(a, b) \cap H_1(a, b)$ is a solution to the boundary value problem (14.1), (14.2). Then,

$$\frac{d}{dx} p(x) \frac{du}{dx} + r(x)u = f(x), \tag{14.11}$$

for almost every $x \in (a, b)$ (see the discussion prior to Example 11.1 for a definition of **almost every**). Multiplying this equality by an arbitrary function $v \in H_0(a, b)$, and integrating over (a, b) , we conclude that

$$\int_a^b \frac{d}{dx} p(x) \frac{du}{dx} v \, dx + \int_a^b r(x)uv \, dx = \int_a^b f(x)v(x) \, dx.$$

On integration by parts in the first term on the left-hand side,

$$\int_a^b \frac{d}{dx} p(x) \frac{du}{dx} v \, dx = p(x) \frac{du}{dx} v \Big|_a^b - \int_a^b p(x) \frac{du}{dx} \frac{dv}{dx} \, dx.$$

Since, by hypothesis, $v(a) = 0$ and $v(b) = 0$, it follows that

$$\int_a^b p(x) \frac{du}{dx} \frac{dv}{dx} \, dx + \int_a^b r(x)uv \, dx = \int_a^b f(x)v(x) \, dx$$

for all $v \in H_0(a, b)$. Thus, we have shown the following result.

Theorem 14.2 *If $u \in H_0(a, b) \cap H_1(a, b)$ is a solution to the boundary value problem (14.1), (14.2), then u is a weak solution to this problem; that is,*

$$(u, v) = (f, v) \quad \forall v \in H_0(a, b). \tag{14.12}$$

The converse implication, namely that any weak solution $u \in H_0(a, b) \cap H_1(a, b)$ of (14.1), (14.2) belongs to $H_0(a, b) \cap H_1(a, b)$ and solves (14.1), (14.2) in the usual (pointwise) sense, is not true in general, unless the weak

the problem by minimising $\mathcal{J}(u)$ over a finite-dimensional subset $S_{\mathbf{g}}^*$ of $H_{\mathbf{g}}(a, b)$, instead.

A simple way of constructing $S_{\mathbf{g}}^*$ is to choose any function $v \in H_{\mathbf{g}}(a, b)$, for example,

$$v(x) = \frac{B}{b} \frac{A}{a} (x - a) + A \tag{14.13}$$

and a finite set of linearly independent functions $v_j, j = 1, \dots, n - 1$, in $H(a, b)$ for $n \geq 2$, and then define

$$S_{\mathbf{g}}^* = \{ v \in H_{\mathbf{g}}(a, b) : v^*(x) = \sum_{j=0}^{n-1} v_j(x), \text{ where } (v_0, \dots, v_{n-1}) \in \mathbb{R}^n \}.$$

We consider the following approximation of problem (RR):

$$(RR)^* \quad \text{find } u^* \in S_{\mathbf{g}}^* \text{ such that } \mathcal{J}(u^*) = \min_{S_{\mathbf{g}}^*} \mathcal{J}(w^*).$$

Our next result is a finite-dimensional analogue of Theorem 14.1.

Theorem 14.4 *A function $u^* \in S_{\mathbf{g}}^*$ minimises $\mathcal{J}(u)$ over $S_{\mathbf{g}}^*$ if, and only if,*

$$(G)^* \quad (\mathcal{J}'(u^*), v^*) = 0, \quad \forall v^* \in S_{\mathbf{g}}^*. \tag{14.14}$$

Here,

$$S^* = \{ v \in H(a, b) : v^*(x) = \sum_{j=0}^{n-1} v_j(x), \text{ where } (v_0, \dots, v_{n-1}) \in \mathbb{R}^n \}.$$

The problem $(G)^*$ can be thought of as an approximation to the Galerkin principle (G), and is therefore referred to as the **Galerkin method**. For a similar reason, $(RR)^*$ is called the Rayleigh–Ritz method, or just **Ritz method**. Thus, in complete analogy with the equivalence of (RR) and (G) formulated in Theorem 14.1, Theorem 14.4 now expresses the equivalence of $(RR)^*$ and $(G)^*$, the approximations to (RR) and (G), respectively. Of course, as in the case of (RR) and (G), the equivalence of $(RR)^*$ and $(G)^*$ relies on the assumption that the bilinear functional $\mathcal{B}(u, v)$ is symmetric. The proof is identical to that of Theorem 14.1, and is left as an exercise.

Theorem 14.4 provides no information about the existence and uniqueness of u^* that minimises $\mathcal{J}(u)$ over $S_{\mathbf{g}}^*$ (or, equivalently, of the existence

and uniqueness of u^* that satisfies (14.14)). This question is settled by our next result.

Theorem 14.5 *There exists a unique function $u^* \in S_{\mathfrak{g}}^*$ that minimises $J(u)$ over $S_{\mathfrak{g}}^*$; this u^* is called the **Ritz approximation** to u . Equivalently, there exists a unique function $u^* \in S_{\mathfrak{g}}^*$ that satisfies (14.14); this u^* is called the **Galerkin approximation** to u . The Ritz and Galerkin approximations to u coincide.*

Proof We shall prove the second of these two equivalent statements: we shall show that there exists a unique $u^* \in S_{\mathfrak{g}}^*$ that satisfies (14.14). The proof of uniqueness of $u^* \in S_{\mathfrak{g}}^*$ is analogous to the proof of Theorem 14.3, with $u, \bar{u}, H_{\mathfrak{g}}(a, b)$ and $H(a, b)$, replaced by $u^*, \bar{u}^*, S_{\mathfrak{g}}^*$ and S^* , respectively. Since $S_{\mathfrak{g}}^*$ is finite-dimensional, the uniqueness of u^* satisfying (14.14) implies its existence. \square

Having shown the existence and uniqueness of u^* minimising $J(u)$ over $S_{\mathfrak{g}}^*$ (or, equivalently, satisfying (14.14)), we adopt the following definition.

Definition 14.4 *The functions $\phi_i, i = 1, 2, \dots, n - 1$, appearing in the definitions of $S_{\mathfrak{g}}^*$ and S^* are called the **Galerkin basis functions**.*

Since any function $v^* \in S^*$ can be represented as a linear combination of the Galerkin basis functions $\phi_i, 1 \leq i \leq n - 1$, it is clear that (14.14) is equivalent to

$$J(u^*) = f, \quad 1 \leq i \leq n - 1. \tag{14.15}$$

As u^* belongs to $S_{\mathfrak{g}}^*$, it can be expressed in terms of ϕ_i and the Galerkin basis functions as

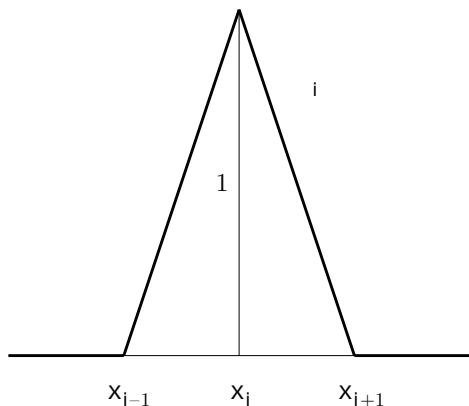
$$u^*(x) = \sum_{i=1}^{n-1} u_i \phi_i(x),$$

where $u_j, j = 1, \dots, n - 1$, are to be determined. On substituting this expansion of u^* into (14.15), we arrive at the following system of simultaneous linear equations:

$$\sum_{j=1}^{n-1} M_{ij} u_j = b_i, \quad 1 \leq i \leq n - 1, \tag{14.16}$$

where

$$\mathbf{M} = \left(\begin{array}{c} \\ \\ \end{array} \right), \quad \mathbf{b} =$$



$$C_{i-1} \phi_{i-1} + C_i \phi_i + C_{i+1} \phi_{i+1}$$

For the finite element method the important property of the basis functions ϕ_i , $1 \leq i \leq n-1$, is that they have *local* support, being nonzero only in one pair of adjacent intervals, $(x_{i-1}, x_i]$ and $[x_i, x_{i+1})$. This means that, in the matrix M ,

$$M_{ij} = 0 \quad \text{if } |i - j| > 1.$$

The matrix M is, therefore, symmetric, positive definite and tridiagonal, and the associated system of linear equations can be solved very efficiently by the methods of Section 3.3, the most efficient algorithm being LU decomposition, without any use of symmetry. The fact that M is positive definite means that no interchanges are necessary.

The function ϕ_i in (14.13), which is included in the definition of S_n^* to ensure that u^* satisfies the boundary conditions at $x = a$ and $x = b$, is then given by

$$\phi_i(x) = A_i \phi_{i-1}(x) + B_i \phi_i(x),$$

which is also piecewise linear; clearly, $\phi_0(a) = A_0$ and $\phi_{n-1}(b) = B_{n-1}$. Here, A_i and B_i are defined by setting, respectively, $i = 0$ and $i = n$ in (14.19) and restricting the resulting functions to the interval $[a, b] = [x_0, x_n]$. In (14.17) we see that the term ϕ_i is nonzero only for $i = 1$ and $i = n - 1$.

Before attempting to solve the system of linear equations we must, of course, first compute the elements of the matrix M , and the quantities on the right-hand side, b_i , $i = 1, \dots, n - 1$; see (14.16) and (14.17). The

matrix elements are obtained from

$$M_{ij} = \int_a^b \left(p(x) \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} + r(x) \phi_i \phi_j \right) dx,$$

with $1 \leq i, j \leq n - 1$. We have written this as the sum of two terms, as the matrix M is often written in this way as the sum of two matrices which, for historical reasons, are often known as the **stiffness matrix** and the **mass matrix**, respectively. The terms M are very simple; in fact in the first integral the derivatives $\frac{d\phi_i}{dx}$ and $\frac{d\phi_j}{dx}$ are piecewise constant functions over $[a, b]$.

It may be possible to compute these integrals analytically, but more generally some form of numerical quadrature will be necessary. It is then easy to show that if we use certain types of quadrature formulae we shall be led to the same system of equations as in the finite difference method of Section 13.5. Consider the particularly simple case where the mesh points are equally spaced, so that $x_j = a + jh$, $j = 0, 1, \dots, n$, $h = (b - a)/n$. If we then approximate the integrals involved in the stiffness matrix by the midpoint rule (see Chapter 10), we obtain

$$\int_a^b p(x) \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} dx = \left(\frac{1}{h} \right) p(x_{ij}) \frac{d\phi_i}{dx} \frac{d\phi_j}{dx} h,$$

where $p_{ij} = p(x_{ij})$, and similarly for the other integrals involved. For the integrals in the mass matrix we use the trapezium rule, and then

$$\int_a^b r(x) \phi_i \phi_j dx = 0,$$

since ϕ_i is zero at x_{j-1} and ϕ_j is zero at x_j . In the same way

$$\int_a^b r(x) \phi_i dx = -hr_i,$$

where $r_i = r(x_i)$, since ϕ_i is zero at one end of the interval and unity at the other. The other part of the integral is, similarly,

$$\int_a^b r(x) \phi_j dx = -hr_j. \tag{14.20}$$

Assuming that $f \in C[a, b]$, approximating the integral on the right-hand side by the trapezium rule in the same way, and putting all the parts together, equation (14.14) now takes the approximate form

$$\frac{p_{ij}}{h} u_{ij} + \frac{p_{ij} + p_{i+1,j}}{h} u_{i+1,j} - \frac{p_i}{h} u_i + hr_i u_i = hf_i,$$

for $i = 1, 2, \dots, n-1$, with the notational convention that $u = A$ and $u = B$, and $f = f(x)$; clearly, this is the same as the finite difference equation (13.19). Of course, had we used a different set of basis functions $\phi_i, 1 \leq i \leq n-1$, or different numerical quadrature rules, the finite element and finite difference methods would have no longer been identical. Indeed, this example is just an illustration of the relation between the two methods; we should normally expect to compute the entries of the matrix M by using some more accurate quadrature method, such as a two-point Gauss formula.

In the next two sections we shall assess the accuracy of the finite element method. Our goal is to quantify the amount of reduction in the error $u - u^*$ as the mesh spacing h is reduced.

14.4 Error analysis of the finite element method

We begin with a fundamental result that underlies the error analysis of finite element methods.

Theorem 14.6 (Céa's Lemma) *Suppose that u is the function that minimises $J(u)$ over $H_g(a, b)$ (or, equivalently, that u satisfies (14.6)), and that u^* is its Galerkin approximation obtained by minimising $J(\cdot)$ over S_g^* (or, equivalently, that u^* satisfies (14.14)). Then,*

$$(u - u^*, v^*) = 0 \quad \forall v^* \in S^*, \quad (14.21)$$

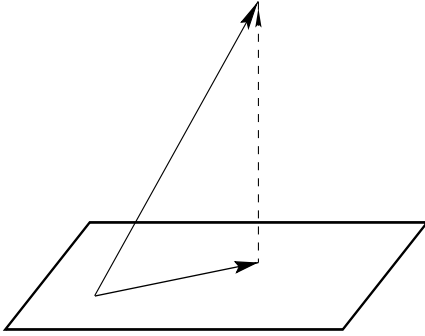
and

$$(u - u^*, u - u^*) = \min_5 (u - v^*, u - v^*). \quad (14.22)$$

The identity (14.21) is referred to as **Galerkin orthogonality**. The terminology stems from the fact that, since the bilinear functional (\cdot, \cdot) is symmetric and $(v, v) > 0$ for all $v \in H(a, b) \setminus \{0\}$, (\cdot, \cdot) is an inner product in the linear space $H(a, b)$. Therefore, by virtue of Definition 9.2, (14.21) means that $u - u^*$ is orthogonal to S^* in $H(a, b)$. A geometrical illustration of Galerkin orthogonality is given in Figure 14.2. Given that u is a fixed element of $H_g(a, b)$, the mapping

$$R^*: u \in H(a, b) \rightarrow u^* \in S^*$$

which assigns a $u^* \in S_g^*$ to $u \in H_g(a, b)$ (where u and u^* are as in Theorem 14.6) is called the **Ritz projector**.



Motivated by the minimisation property (14.22), we define the **energy norm** $\| \cdot \|_A$ on $H_0^1(a, b)$ via

$$\|v\|_A = \left((v, v) \right)^{1/2}. \tag{14.23}$$

Under our hypotheses on p and q , it is easy to see that $\| \cdot \|_A$ satisfies all axioms of norm (see Chapter 2). The result we have just proved shows that u^* is the *best approximation* from S_h^* to the true solution $u \in H_0^1(a, b)$ of our problem, when we measure the error of the approximation in the energy norm:

$$\|u - u^*\|_A = \min_{v \in S_h^*} \|u - v\|_A. \tag{14.24}$$

A particularly relevant question is how the error $\|u - u^*\|_A$ depends on the spacing h of the subdivision of the computational domain $[a, b]$. We can obtain a bound on the error $\|u - u^*\|_A$, measured in the energy norm, by choosing a particular function $v^* \in S_h^*$ in (14.24) whose closeness to u is easy to assess. For this purpose, we introduce the **finite element interpolant** $u^h \in S_h^*$ of $u \in H_0^1(a, b)$ by

$$u^h(x) = \sum_{j=0}^n u(x_j) \phi_j(x), \quad x \in [a, b].$$

Clearly,

$$u^h(x_j) = u(x_j), \quad j = 0, 1, \dots, n,$$

which justifies our use of the word *interpolant*.

We then deduce from (14.24) that

$$\|u - u^*\|_A \leq \|u - u^h\|_A; \tag{14.25}$$

hence, in order to quantify $\|u - u^*\|_A$, we only need to estimate the size of $\|u - u^h\|_A$. This leads us to the next theorem.

Theorem 14.7 *Suppose that $u \in H_0^1(a, b) \cap H_0^2(a, b)$ and let u^h be the finite element interpolant of u from S_h^* defined above; then, the following error bounds hold:*

$$\|u - u^h\|_A \leq \frac{h}{2} \|u''\|_A, \tag{14.26}$$

$$\|u - u^h\|_{L^\infty} \leq \frac{h^2}{8} \|u''\|_A,$$

for $i = 1, 2, \dots, n$, where $h_i = x_i - x_{i-1}$.

Proof Consider an element $[x_{i-1}, x_i]$, $1 \leq i \leq n$, and define $\phi_i(x) = u(x_i) - u(x)$ for $x \in [x_{i-1}, x_i]$. Then, $\phi_i(x_i) = 0$ and $\phi_i(x_{i-1}) = u(x_{i-1}) - u(x_i) = -\Delta u_i$. Therefore ϕ_i can be expanded into a convergent Fourier sine-series,

$$\phi_i(x) = \sum_{k=1}^{\infty} a_k \sin\left(\frac{k\pi(x - x_{i-1})}{h}\right)$$

Now, substituting the bounds from Theorem 14.7 into the definition of the norm $\|u - u_A\|_A$, we arrive at the following estimate of the interpolation error in the energy norm.

Corollary 14.1 *Suppose that $u \in H^1(a, b) \cap H_{\&}(a, b)$. Then,*

$$\|u - u_A\|_A \leq \frac{h}{2} P + \frac{h}{2} R \|u\|_{\#},$$

where $P = \max_{x \in [a, b]} p(x)$ and $R = \max_{x \in [a, b]} r(x)$.

Proof Let us observe that

$$\begin{aligned} \|v\|_A^2 &= \int_a^b (v, v) \\ &= \int_a^b p(x) v(x)^2 + r(x) v(x)^2 \, dx \\ &= \int_a^b p(x) v(x)^2 + r(x) v(x)^2 \, dx \\ &= \frac{3}{4} P \|v\|_{\#}^2 + \frac{1}{4} R \|v\|_{\#}^2. \end{aligned}$$

On letting $v = u - u_A$ and applying the preceding theorem on the right-hand side of the last inequality, with v and v replaced by $u - u_A$ and $u - u_A$, respectively, the result follows. \square

Inserting this estimate into (14.25) leads to the desired bound on the error between the analytical solution u and its finite element approximation u_A in the energy norm.

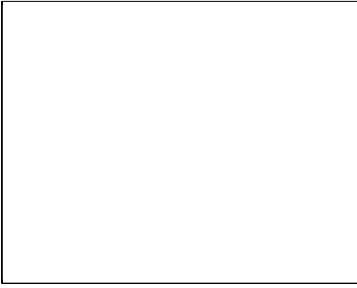
Corollary 14.2 *Suppose that $u \in H^1(a, b) \cap H_{\&}(a, b)$. Then,*

$$\|u - u_A\|_A \leq \frac{h}{2} P + \frac{h}{2} R \|u\|_{\#},$$

where $P = \max_{x \in [a, b]} p(x)$ and $R = \max_{x \in [a, b]} r(x)$. Further,

$$\|u - u_A\|_A \leq \frac{h}{2} P + \frac{h}{2} R \|u\|_{\#}, \tag{14.26}$$

where $P = \max_{x \in [a, b]} p(x)$, $R = \max_{x \in [a, b]} r(x)$, and $h = \max_{x \in [a, b]} h$.



We see from Figure 14.3 that, as the spacing h of the subdivision is reduced, the finite element solution u^* approximates the analytical solution $u(x) = \sin(x)$ with increasing accuracy. Indeed, the results corresponding to $n = 2$ and $n = 4$ in Figure 14.3 indicate that as the number of intervals in the subdivision is doubled (*i.e.*, h is halved), the maximum error between $u(x)$ and $u^*(x)$ is reduced by a factor of about 4. This reduction in the error cannot be explained by Corollary 14.2 which merely implies that halving h should lead to a reduction in $\|u - u^*\|_A$ by a factor no less than 2. If you would like to learn more about the source of the observed enhancement of accuracy, consult Exercise 5 at the end of the chapter.

14.5 error analysis by duality

The bound on the error between the analytical solution u and its finite element approximation u^* formulated in Corollary 14.2 shows that, in the limit of $h \rightarrow 0$, the error $\|u - u^*\|_A$ will tend to zero as $h \rightarrow 0$. This is a useful result from the theoretical point of view: it reassures us that the unknown analytical solution may be approximated arbitrarily well by making h sufficiently small. On the other hand, asymptotic error bounds of this kind are not particularly helpful for the purpose of precisely quantifying the size of the error between u and u^* for a given, *fixed*, mesh size $h > 0$: as u is unknown, it is difficult to tell just how large the right-hand side of (14.26) really is.

The aim of the present section is, therefore, to derive a computable bound on the error, and to demonstrate how such a bound may be implemented into an adaptive mesh-refinement algorithm, capable of reducing the error $\|u - u^*\|_A$ below a certain prescribed tolerance in an automated manner, without human intervention. The approach is based on seeking a bound on $\|u - u^*\|_A$ in terms of the computed solution u^* rather than in terms of norms of the unknown analytical solution u . A bound on the error in terms of u^* is referred to as an *error bound*, due to the fact that it becomes *computable* only *after* the numerical solution u^* has been obtained.

In order to illuminate the key ideas while avoiding technical difficulties, we shall consider the two-point boundary value problem

$$(p(x)u') + q(x)u + r(x)u = f(x), \quad a < x < b, \quad (14.28)$$

$$u(a) = A, \quad u(b) = B, \quad (14.29)$$

where $p, q \in C[a, b]$, $r \in C[a, b]$ and $f \in L(a, b)$. We shall assume, as

at the beginning of the chapter, that $p(x) = c > 0$, $x \in [a, b]$; however, instead of supposing that $r(x) = 0$, we shall now demand that

$$r(x) = \frac{1}{2}q(x) = c, \quad x \in [a, b], \quad (14.30)$$

where c is assumed to be a positive constant.

Letting

$$(w, v) = \int_a^b [p(x)w + r(x)v] dx$$

consider the auxiliary boundary value problem

$$(p(x)z)' - (q(x)z) + r(x)z = (u - u^*)(x), \quad a < x < b, \quad (14.33)$$

$$z(a) = 0, \quad z(b) = 0, \quad (14.34)$$

called the **dual problem** (or adjoint problem).

We begin our error analysis by noting that the definition of the dual problem and straightforward integration by parts yield (recalling that $(u - u^*)(a) = 0, (u - u^*)(b) = 0$)

$$\begin{aligned} (u - u^*)' - (qz) + rz &= (u - u^*)' - (qz) + rz \\ &= (u - u^*)' - (qz) + rz \\ &= ((u - u^*)', z). \end{aligned}$$

On the other hand, (14.31) and (14.32) imply the Galerkin orthogonality property

$$((u - u^*)', z^*) = 0 \quad \forall z^* \in S^*.$$

In particular, by choosing

$$z^* = \Pi_h z \in S^*,$$

the continuous piecewise linear interpolant of the function $z \in H^1(a, b)$, associated with the subdivision $a = x_0 < x_1 < \dots < x_{n-1} < x_n = b$, we have that

$$((u - u^*)', \Pi_h z) = 0.$$

Thus,

$$\begin{aligned} (u - u^*)' - (qz) + rz &= ((u - u^*)', \Pi_h z) \\ &= ((u - u^*)', \Pi_h z) - ((u^* - u^*)', \Pi_h z) \\ &= ((u - u^*)', \Pi_h z) - ((u^* - u^*)', \Pi_h z), \end{aligned} \quad (14.35)$$

where the last transition follows from (14.31) with $v = \Pi_h z$.

We observe that the right-hand side no longer involves the unknown analytical solution u . Furthermore,

$$\begin{aligned} ((u^* - u^*)', \Pi_h z) &= \int_a^b p(x)(u^* - u^*)(x)(z - \Pi_h z)(x) dx \\ &+ \int_a^b q(x)(u^* - u^*)(x)(z - \Pi_h z)(x) dx \\ &+ \int_a^b r(x)u^*(x)(z - \Pi_h z)(x) dx. \end{aligned}$$

Integrating by parts in each of the n integrals in the first sum on the right-hand side, noting that $(z_{i+1}^*, z_i^*)(x) = 0$, $i = 0, \dots, n$, we deduce that

$$\begin{aligned} & \left((u^*, z_{i+1}^* - z_i^*) \right) \\ &= \int_{x_{i-1}}^{x_i} (p(x)(u^*)' + q(x)(u^*) + r(x)u^* - (z_{i+1}^* - z_i^*)'(x)) dx. \end{aligned}$$

Furthermore,

$$\int_{x_{i-1}}^{x_i} (z_{i+1}^* - z_i^*)'(x) dx = \int_{x_{i-1}}^{x_i} f(x) (z_{i+1}^* - z_i^*)(x) dx.$$

Substituting these two identities into (14.35), we deduce that

$$\|u - u^*\|_{\#} = \int_{x_{i-1}}^{x_i} R(u^*)(x) (z_{i+1}^* - z_i^*)(x) dx, \quad (14.36)$$

where, for $1 \leq i \leq n$, and $x \in (x_{i-1}, x_i)$,

$$R(u^*)(x) = f(x) - (p(x)(u^*)' + q(x)(u^*) + r(x)u^*).$$

The function $R(u^*)$ is called **the finite element residual**; it measures the extent to which u^* fails to satisfy the differential equation

$$(p(x)u)' + q(x)u + r(x)u = f(x)$$

on the union of the intervals (x_{i-1}, x_i) , $i = 1, \dots, n$. Now, applying the Cauchy–Schwarz inequality on the right-hand side of (14.36) yields

$$\|u - u^*\|_{\#} \leq \|R(u^*)\|_{\#} \|z_{i+1}^* - z_i^*\|_{\#}.$$

Recalling from Theorem 14.7 that

$$\|z_{i+1}^* - z_i^*\|_{\#} \leq \frac{h}{2} \|z_i^*\|_{\#}, \quad i = 1, 2, \dots, n,$$

we deduce that

$$\|u - u^*\|_{\#} \leq \frac{1}{2} h \|R(u^*)\|_{\#} \|z_i^*\|_{\#},$$

and consequently, using the Cauchy–Schwarz inequality for finite sums,

$$\|a\| \|b\| \leq \|a\| \|b\|$$

with

a

Integrating by parts, again, in the second term on the right gives

$$\begin{aligned}
 (pz) - (qz) + rz, z! &= c z_{\#} - \frac{1}{2} \int q(x)[z'(x)] dx \\
 &+ \int r(x)[z(x)] dx.
 \end{aligned}$$

Hence, from (14.39),

$$c z_{\#} + \int r(x) \frac{1}{2} q(x) [z(x)] dx = u - u^*, z!,$$

and thereby, noting (14.30) and using the Cauchy–Schwarz inequality on the right-hand side,

$$\min c, c z_{\#} + z_{\#} = u - u^*, z! \leq u - u^*, z_{\#}. \tag{14.40}$$

Therefore, also

$$\min c, c z_{\#} = u - u^*, z_{\#},$$

which means that

$$z_{\#} + z_{\#} = z_{\#} \frac{1}{\min c, c} u - u^*, z_{\#}. \tag{14.41}$$

Now we substitute (14.41) into (14.38) to deduce that

$$z_{\#} = K u - u^*, z_{\#}, \tag{14.42}$$

where

$$K = \frac{1}{c} \left(1 + \frac{1}{\min c, c} \int p + q + r \right) q.$$

□

It is important to observe here that K involves only known quantities: the coefficients in the differential equation under consideration. Therefore K can be computed, or at least bounded above, without difficulties. On inserting (14.42) into (14.37), we arrive at our final result, the computable *a posteriori* error bound,

$$u - u^*_{\#} = K h R(u^*)_{\#}, \tag{14.43}$$

where $K = K/\dots$.

Next we shall describe the construction of an adaptive mesh refinement algorithm based on the *a posteriori* error bound (14.43).

Suppose that **012**

- (4 If not, then halve those elements $[x_{i-1}, x_i]$ in \mathcal{T}_m , with i in the set $\{1, 2, \dots, n\}$, for which

$$h_i = x_i - x_{i-1} > \frac{1}{n} \frac{012}{K}, \quad (14.48)$$

denote by \mathcal{T}_{m+1} the resulting subdivision of $[a, b]$ with n elements $[x_{i-1}, x_i]$ of respective lengths

$$h_i = x_i - x_{i-1}, \quad i = 1, \dots, n,$$

and consider the associated finite element space S_m^* of dimension $n - 1$;

- (! Compute the finite element approximation $u^* \in S_m^*$, increase m by 1 and return to (4,

The inequality (14.47) is called the **stopping criterion** for the mesh adaptation algorithm, and (14.48) is referred to as the **refinement criterion**. According to the *a posteriori* error bound (14.43), when the adaptive algorithm terminates, the error $\|u - u^*\|_{\infty}$ is guaranteed not to exceed the prescribed tolerance 012 .

We conclude the body of this chapter with a numerical experiment which illustrates the performance of the adaptive algorithm.

Example 14.1 *Let us consider the second-order ordinary differential equation*

$$(p(x)u)' + q(x)u + r(x)u = f(x), \quad x \in (0, 1), \quad (14.49)$$

subject to the boundary conditions

$$u(0) = 0, \quad u(1) = 0. \quad (14.50)$$

Suppose, for example, that

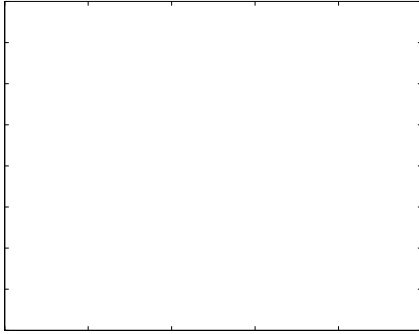
$$p(x) = 1, \quad q(x) = 20, \quad r(x) = 10 \quad \text{and} \quad f(x) = 1.$$

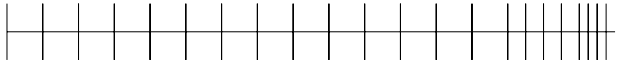
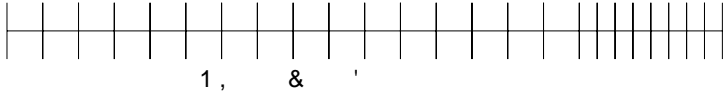
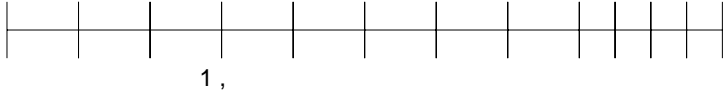
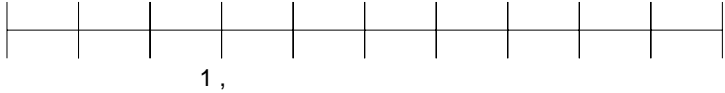
In this case, the analytical solution, u , can be expressed in closed form:

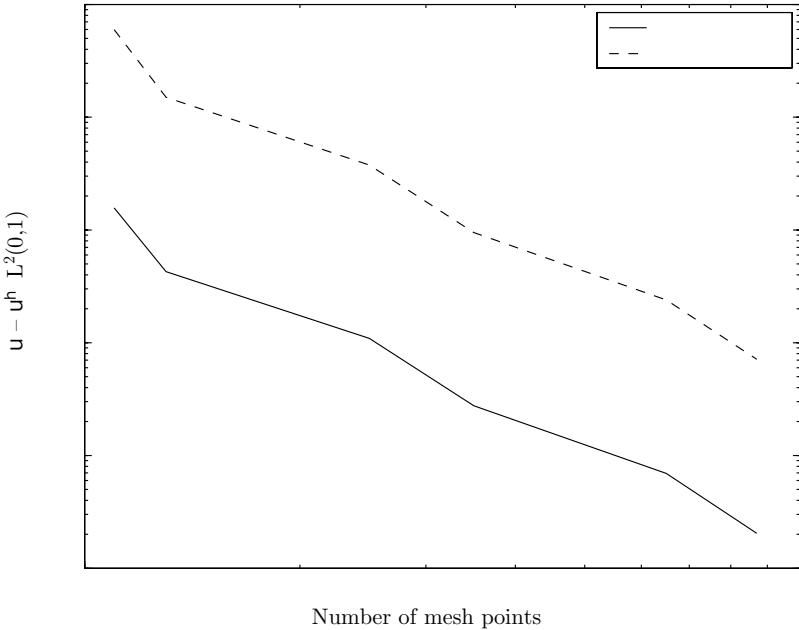
$$u(x) = C_1 e^{\alpha x} + C_2 e^{\beta x} + \frac{1}{10},$$

where α and β are the two roots of the characteristic polynomial of the differential equation, $\lambda^2 + 20\lambda + 10 = 0$, i.e.,

$$\alpha = 10 + \sqrt{110}, \quad \beta = 10 - \sqrt{110},$$







'% 8 L (0 1) * +
 #%&# 4', ,B * +

, / - " / "\$ " , Introduction to adaptive methods for differential equations, in *Acta Numerica* **4** (A. Iserles, ed.), Cambridge University Press, Cambridge, 105–158, 1995.

% ! - % % , An optimal control approach to a-posteriori error estimation in finite element methods, in *Acta Numerica* **10** (A. Iserles, ed.), Cambridge University Press, Cambridge, 1–102, 2001.

+ ! / &5 , Adjoint methods for PDEs: superconvergence and adaptivity by duality, in *Acta Numerica* **11** (A. Iserles, ed.), Cambridge University Press, Cambridge, 145–236, 2002.

A detailed and general survey of the subject of *a posteriori* error estimation can be found in

+) (, *A posteriori Error Estimation in Finite Element Analysis*, John Wiley & Sons, New York, 2000.

In this chapter we were concerned with the *a priori* error analysis of the piecewise linear finite element method in the energy norm, and its *a posteriori* error analysis in the L^2 norm. Using similar techniques, one can establish an *a priori* error bound in the L^2 norm and an *a posteriori* error bound in the energy norm. For extensions of the theory considered here to higher-order piecewise polynomial finite element approximations and generalisations to partial differential equations, the reader is referred to the books listed above.

Exercises

14.1 Given that (a, b) is an open interval of the real line, let

$$H_1(a, b) = \{v \in H^1(a, b) : v(a) = 0\}.$$

(i) By writing

$$v(x) = \int_a^x v'(t) dt,$$

for $v \in H_1(a, b)$ and $x \in [a, b]$, show the following (**Poincaré–Friedrichs**) inequality:

$$\|v\|_{L^2(a,b)} \leq \frac{1}{2}(b-a)^{1/2} \|v'\|_{L^2(a,b)}.$$

(ii) By writing

$$[v(x)]^2 = \int_a^x 2v(t)v'(t) dt = 2 \int_a^x v(t)v'(t) dt$$

for $v \in H_1(a, b)$ and $x \in [a, b]$, show the following (**Agmon’s**) inequality:

$$\max_{x \in [a,b]} |v(x)| \leq 2 \|v'\|_{L^2(a,b)}.$$

14.2 Given that $f \in L^2(0, 1)$, state the weak formulation of each of the following boundary value problems on the interval $(0, 1)$:

- (a) $-u'' + u = f(x), u(0) = 0, u(1) = 0$;
- (b) $-u'' + u = f(x), u(0) = 0, u(1) = 1$;
- (c) $-u'' + u = f(x), u(0) = 0, u(1) + u'(1) = 2$.

In each case, show that there exists at most one weak solution.

14.3 Give a proof of Theorem 14.4.

14.4 Prove Corollary 14.2.

14.5 Consider the boundary value problem

$$p u'' + r u' = f(x), \quad u(0) = 0, \quad u(1) = 0,$$

on the interval $[0, 1]$, where p and r are positive constants and $f \in C[0, 1]$. Using equally spaced points

$$x_i = ih, \quad i = 0, 1, \dots, n, \quad \text{with } h = 1/n, \quad n \geq 2,$$

and the standard piecewise linear finite element basis functions (hat functions) $\phi_i, i = 1, 2, \dots, n-1$, show that the finite element equations for $u = u^*(x)$ become

$$p(u_{i-1} - 2u_i + u_{i+1})/h + r(u_{i-1} + 4u_i + u_{i+1})/6 = \frac{1}{h} f_i, \quad i = 1, 2, \dots, n-1,$$

with $u_0 = 0$ and $u_n = 0$. By expanding in Taylor series, show that

$$\frac{1}{h} f_i = f(x_i) + \frac{h^2}{6} f''(x_i) + O(h^4).$$

Interpreting this set of difference equations as a finite difference approximation to the boundary value problem, as in Chapter 13, show that the corresponding truncation error T satisfies

$$T_i = -\frac{h^2}{6} r u''(x_i) + O(h^4), \quad i = 1, \dots, n-1,$$

and use the method of Exercise 13.2 to show that

$$\max |u(x) - u^*(x)| \leq Mh^2,$$

where M is a positive constant.

14.6 In the notation of Exercise 5 suppose that all the integrals involved in the calculation are approximated by the trapezium rule. Show that the system of equations becomes identical to that obtained from the central difference approximation in Chapter 13, and deduce that

$$\max |u(x) - u^*(x)| \leq Mh^2,$$

where M is a positive constant.

14.7 Consider the differential equation

$$(p(x)u')' + r(x)u = f(x), \quad a < x < b,$$

with p, r and f as at the beginning of the chapter, subject to the boundary conditions

$$p(a)u'(a) + u(a) = A, \quad p(b)u'(b) + u(b) = B,$$

where α and β are positive real numbers, and A and B are real numbers. Show that the weak formulation of the boundary value problem is

find $u \in H^1(a, b)$ such that $(u, v) = (f, v)$ for all $v \in H^1(a, b)$, where

$$(u, v) = \int_a^b [\alpha(x)u'(x)v'(x) + \gamma(x)u(x)v(x)]dx + \alpha(a)u(a)v(a) + \alpha(b)u(b)v(b),$$

and

$$(f, v) = \int_a^b f(x)v(x)dx + Av(a) + Bv(b).$$

Construct a finite element approximation of the boundary value problem based on this weak formulation using piecewise linear finite element basis functions on the subdivision

$$a = x_0 < x_1 < \dots < x_n = b$$

of the interval $[a, b]$. Show that the finite element method gives rise to a set of $n + 1$ simultaneous linear equations with $n + 1$ unknowns $u_i = u^*(x_i)$, $i = 0, 1, \dots, n$. Show that this linear system has a unique solution.

Comment on the structure of the matrix M of the linear system: (a) Is M symmetric? (b) Is M positive definite? (c) Is M tridiagonal?

14.8 Given that α is a nonnegative real number, consider the differential equation

$$-u'' + \alpha u = f(x) \quad \text{for } x \in (0, 1),$$

subject to the boundary conditions

$$u(0) = 0, \quad u(1) + u'(1) = 0.$$

State the weak formulation of the problem. Using continuous piecewise linear basis functions on a uniform subdivision of $[0, 1]$ into elements of size $h = 1/n$, $n \geq 2$, write down the finite element approximation to this problem and show that this has a unique solution u^* . Expand u^* in terms of the standard piecewise linear finite element basis functions (hat functions) ϕ_i ,

$i = 1, 2, \dots, n$, by writing

$$u^*(x) = \sum_{i=1}^n U_i(x)$$

to obtain a system of linear equations for the vector of unknowns (U_1, \dots, U_n) .

Suppose that $\alpha = 0$, $f(x) = 1$ and $h = 1/3$. Solve the resulting system of linear equations and compare the corresponding numerical solution $u^*(x)$ with the exact solution $u(x)$ of the boundary value problem.

14.9 Consider the differential equation

$$(p(x)u') + r(x)u = f(x), \quad x \in (0, 1),$$

subject to the boundary conditions $u(0) = 0$, $u(1) = 0$, where $p(x) > 0$, $r(x) \geq 0$ for all x in the closed interval $[0, 1]$, with $p \in C^1[0, 1]$, $r \in C[0, 1]$ and $f \in L^1(0, 1)$. Given that u^* denotes the continuous piecewise linear finite element approximation to u on a uniform subdivision of $[0, 1]$ into elements of size $h = 1/n$, $n \geq 2$, show that

$$\|u - u^*\|_{\infty} \leq C h \|u\|_{\infty},$$

where C is a positive constant that you should specify. Show further that there exists a positive constant C such that

$$\|u - u^*\|_{\infty} \leq C h \|f\|_{\infty}.$$

Calculate the right-hand sides in these inequalities in the case when

$$p(x) = 1, \quad r(x) = 0, \quad f(x) = 1,$$

for $x \in [0, 1]$, and $h = 10^{-2}$.

14.10 Consider the two-point boundary value problem

$$-u'' + u = f(x), \quad x \in (0, 1), \quad u(0) = 0, \quad u(1) = 0,$$

with $f \in C[0, 1]$. State the piecewise linear finite element approximation to this problem on a nonuniform subdivision

$$0 = x_0 < x_1 < \dots < x_n = 1, \quad n \geq 2,$$

with $h_i = x_i - x_{i-1}$, assuming that, for a continuous piecewise

linear function v^* ,

!

$$\int_{x_{i-1}}^{x_i} f(x)v^*(x)dx$$

has been approximated by applying the trapezium rule on each element $[x_{i-1}, x_i]$.

Verify that the following *a posteriori* bound holds for the error between u and its finite element approximation u^* :

$$\|u - u^*\|_{\infty} \leq K \max_i h_i \|R(u^*)\|_{\infty}$$

$$+ K \max_i h_i \left(\max_{x \in [x_{i-1}, x_i]} |f''(x)| + 4 \max_{x \in [x_{i-1}, x_i]} |f'(x)| \right),$$

where $R(u^*) = f(x) - ((u^*)'(x) + u^*(x))$ for $x \in [x_{i-1}, x_i]$, $i = 1, \dots, n$, and K_1, K_2 are constants which you should specify.

How would you use this bound to compute u to within a specified tolerance 012?

Appendix A

An overview of results from real analysis

In this Appendix we gather a number of results from real analysis which are assumed at various places in the text. Some of these will be familiar from any course on the subject, and no proofs are given; a small number may be less familiar, and we give proofs of these for completeness.

Theorem A.1 (The Intermediate Value Theorem) *Suppose that f is a real-valued function, defined and continuous on the closed interval $[a, b]$ of \mathbb{R} . Then, f is a bounded function on the interval $[a, b]$ and, if y is any number such that*

$$\inf_{x \in [a, b]} f(x) \leq y \leq \sup_{x \in [a, b]} f(x),$$

then there is a number $c \in [a, b]$ such that $f(c) = y$. In particular, the infimum and the supremum of f are achieved, and can be replaced by $\min_{x \in [a, b]} f(x)$ and $\max_{x \in [a, b]} f(x)$, respectively.

The next result, known as Rolle's Theorem, was published in an obscure book in 1691 by the French mathematician Michel Rolle (1652–1719) who invented the notation $\sqrt[n]{x}$ for the n th root of x .

Theorem A.2 (Rolle's Theorem) *Suppose that f is a real-valued function, defined and continuous on the closed interval $[a, b]$ of \mathbb{R} , differentiable in the open interval (a, b) , and such that $f(a) = f(b)$. Then, there exists a number $c \in (a, b)$ such that $f'(c) = 0$.*

It is often important in our applications that the point $c \in (a, b)$, i.e., $a < c < b$. For instance it may happen that $f'(a) = f'(b) = 0$, as well as $f(a) = f(b)$; Theorem A.2 then states that, in addition to the endpoints

of the interval $[a, b]$, there is also an interior point $c \in (a, b)$ at which the derivative vanishes.

Theorem A.3 (The Mean Value Theorem) *Suppose that f is a real-valued function, defined and continuous on the closed interval $[a, b]$ of \mathbb{R} , and f is differentiable in the open interval (a, b) . Then, there exists a number $c \in (a, b)$ such that*

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

Theorem A.4 (Taylor's Theorem) *Suppose that n is a nonnegative integer, and f is a real-valued function, defined and continuous on the closed interval $[a, b]$ of \mathbb{R} , such that the derivatives of f of order up to and including n are defined and continuous on the closed interval $[a, b]$. Suppose further that f is differentiable on the open interval (a, b) . Then, for each value of x in $[a, b]$, there exists a number $c = c(x)$ in the open interval (a, b) such that*

$$\begin{aligned} f(x) = & f(a) + (x - a)f'(a) + \frac{(x - a)^2}{2!}f''(a) \\ & + \frac{(x - a)^3}{3!}f'''(a) + \cdots \\ & + \frac{(x - a)^n}{n!}f^{(n)}(a) \\ & + \frac{(x - a)^{n+1}}{(n+1)!}f^{(n+1)}(c). \end{aligned}$$

Theorem A.5 (Taylor's Theorem with integral remainder) *Let n be a nonnegative integer and suppose that f is a real-valued function, defined and continuous on the closed interval $[a, b]$ of \mathbb{R} , such that the derivatives of f of order up to and including n are defined and continuous on $[a, b]$, f is differentiable on the open interval (a, b) , and $f^{(n+1)}$ is integrable on (a, b) . Then, for each $x \in [a, b]$,*

$$\begin{aligned} f(x) = & f(a) + (x - a)f'(a) + \frac{(x - a)^2}{2!}f''(a) \\ & + \frac{(x - a)^3}{3!}f'''(a) + \cdots \\ & + \frac{(x - a)^n}{n!}f^{(n)}(a) \\ & + \frac{(x - t)^n}{n!}f^{(n+1)}(t)dt. \end{aligned}$$

Proof As this version of the theorem may be rather less familiar we include a proof.

The theorem is trivially true for $n = 0$. Suppose that the theorem is true for some nonnegative integer, say $n = k$. Then, provided that f is differentiable on (a, b) and $f^{(k+1)}$ is integrable on (a, b) , integration

by parts shows that

$$\int_a^x \frac{(x-t)^k}{(k+1)!} f(t) dt = \frac{(x-a)^k}{(k+1)!} f(a) + \int_a^x \frac{(x-t)^{k-1}}{k!} f(t) dt;$$

use of the theorem when $n = k$ now shows that it is also true for $n = k+1$. The proof by induction is then complete. □

Theorem A.6 (The Integral Mean Value Theorem) *Suppose that f is a real-valued function, defined and continuous on a closed interval $[a, b]$ of \mathbb{R} , and let g be a function, defined, nonnegative and integrable on (a, b) . Then, there exists a number $\xi \in (a, b)$ such that*

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

Proof Since f is continuous on $[a, b]$, it is bounded on $[a, b]$, say

$$m \leq f(x) \leq M, \quad x \in [a, b].$$

Then, as $g(x) \geq 0$ for all $x \in (a, b)$, we have that

$$mg(x) \leq f(x)g(x) \leq Mg(x), \quad x \in (a, b).$$

Integrating these inequalities gives

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

If $\int_a^b g(x) dx = 0$, then the result trivially follows. If, on the other hand, $\int_a^b g(x) dx > 0$, then

$$m \leq \frac{\int_a^b f(x)g(x) dx}{\int_a^b g(x) dx} \leq M.$$

The existence of the required value of $\xi \in (a, b)$ now follows from the Intermediate Value Theorem. □

Theorem A.6 obviously also holds provided that $g(x) \leq 0$ on (a, b) ; it is only important that g has constant sign on (a, b) . Note also that we do not require that g is continuous, only that it is integrable. For example, Theorem A.6 will hold if f is a continuous function defined on $[0, 1]$ and $g(x) = x^{-1/2}$, $x \in (0, 1)$.

Theorem A.7 (Taylor's Theorem for several variables) Suppose that f is a real-valued function of n real variables, $n \geq 1$, such that f and all of its partial derivatives up to and including order $k + 1$ are defined, continuous and bounded in a neighbourhood of the point \mathbf{a} in \mathbb{R}^n . Let A denote an upper bound on the absolute values of all the derivatives of order $k + 1$ in this neighbourhood. Then

$$f(\mathbf{a} + \mathbf{h}) = f(\mathbf{a}) + \frac{U(\mathbf{h})}{r!} + E, \quad (1)$$

where

$$U(\mathbf{h}) = \frac{\partial^r f}{\partial x_1 \dots \partial x_r}(\mathbf{a}) + \dots + \frac{\partial^r f}{\partial x_r \dots \partial x_1}(\mathbf{a}), \quad r = 1, \dots, k,$$

and

$$|E| \leq \frac{A}{(k+1)!} \|\mathbf{h}\|^{k+1}.$$

Proof The proof involves the application of Theorem A.4, Taylor's Theorem, to the function of one variable

$$f(\mathbf{a} + t\mathbf{h})$$

to give a series expansion for (1). Then, the expressions for the derivatives of $f(\mathbf{a} + t\mathbf{h})$ in terms of the partial derivatives of f , via the chain rule, yield the required result; n is the number of partial derivatives of order $k + 1$ for a function of n variables. \square

Appendix B

WWW-resources

The book would not be complete without some mention of numerical analysis software and software repositories on the World Wide Web.

An excellent source of mathematical software is the Netlib Repository on the website

A detailed classified list of the available mathematical software libraries can be viewed by clicking on the - button on this webpage. It is also possible to search the repository for a specific piece of software.

Another useful resource is the website of the *ACM Transactions on Mathematical Software* (TOMS) at

5

The site maintains a well-organised repository, including a range of freely available packages for both numerical and symbolical computations, as well as a number of helpful links to various software vendors. The latter include the developers of Maple (a software for symbolical and numerical computations, scientific visualisation and programming), the makers of Mathematica (a software system for symbolical, numerical and graphical computations), the Numerical Algorithms Group (NAG), MathWorks, Inc., the developers of Matlab (a technical computing environment for high-performance numerical computation and visualisation), and many others. Most of the numerical experiments included in the book were performed by using either Matlab or Maple.

Concerning the history of mathematics, we refer to the Mac Tutor history of mathematics website at St Andrews University in Scotland:

A more recent site, dedicated specifically to the history of approximation theory, resides on

Bibliography

C, : F 1 7C *\$)+ / *E
/ K >+
C, L9 / -/ M L3 *\$%)+
F *C
@ / K >+
C F ? L5 * +
*L M BN1 / K >+
C 4F F A <FF < <E *\$\$(+
! " # \$% &
8 \$ *17CF @ +
C! ? *\$\$%+ ' (*8, 4 B @
8, +
2 * &+)
) * \$ +,
- ' *) * ,
(B @ 8, + *8, 4 "
2 3 > 85 *\$%+ / ! + +
&|)'
2 > < < < * + C "
I C 7 *8, 4 B @ 8, +
2 1/ *\$ O&+ EP P Q M P
,, P) + + / I
2AR- @ *\$\$%+ 0 & " *17CF
@ +
2B 51 <, - *\$\$+ 0 * *1 3 "
+
2 E * + *8, 4 B @ 8, +
2 1 1 3 < * +
1 - *1 / K >+
2 L8 *\$)+ # \$% &

+ 1 *2/ 0 *L M B N 1

8 +

8 , L "3 *\$\$\$+ * "

*1 / K >+

8 B - M *\$%+ ' - *F "9

/ K >+

8 @ *\$\$\$+ ' 0 * #

*8 , 4 B @ 8 , +

8 < *\$&+ T , G "

, , + + + + I&

8 *) (+ ' 3 4 * *5 &

*8 : Q 8 N8 @ , Q +

E G *\$%&+ C , B ,

')I&

E @L < , : @ *\$'+ ' * 1 "

- *C @ ? 3+

E 2 8 * + " < - *1

/ K >+

E L"@ 1 , F * + F /

+) +)I \$

E >> J T L *\$'+ 1 *2/ %

\$% & */ "9 C +

E 1C *\$\$\$+ ' *C I

M B < FC+

E : @ *\$\$\$+ *8 , 4 B @ 8 "

, +

EB @1 *\$'+ C ! G

B ! A +

+ + I\$

- 9 *\$+ 6) 8 F "

C *C @ 3 +

- > J - E 9 , @ L 8 *\$\$(+ 7

H G (I (

C 7 *8 , 4 B @ 8 , +

*\$)+ # * " * @ @ ?! /

K >+

A 1 L *\$\$\$+ R, 1C *\$ I\$&&+

+ + * + + + (" \$

8 * \$+

+ * @ 79 2 9 , +

8 * \$+ F !

8 %&I\$% *E R %&+

M *\$\$\$+ ? B C

C 7 *8 , 4 B @ 8 "

, +

M *\$\$\$)+ 7 ' *2 >R

2 FC+

M , 9 ? *\$\$\$+)
 # * " *2 >R 2 +
 F2 1R - * + CA @E- "
 B,B B '(I&% C 7
 *8 , 4 B@ 8 , +
 9 *\$))+ - *
) *1 / K >+
 , 9 T 3 8 *\$\$%+ -) 5 -
 *L 9 > 4 B@ 2 +
 9 C *\$ + E F > > C R
 > + + \$!&
 9 - /U 1@ M *\$\$&+ (*# \$%
 & '7 %" 1 - *1 2 +
 9 - M *\$\$\$+ (*# \$% & '7
 % \$% * " *1 2 +
 9 C *\$\$\$+ 0 * *8 , 4 B@ 8 "
 , +
 9 @ *\$%+ \$! # \$% &
 *L M BN1 / K >+
 9 /L *\$\$%+ * *17CF
 @ +
 9 , 2 *\$(%+ ' *F 19
 / K >+
 9 B L 8 *\$(+ 1 ! +) +
 &&I'&
 9 <C L 8 < *\$\$\$+ - *8 , 4 B
 @ 8 , +
 9 C1 *\$%'+
 *2 / K >+
 9 C1 *\$)+ *
 & *F 19 / K >+
 7 C *\$\$\$+) \$%
 & *8 , 4 B@ 8 , +
 L > FC 5 , L *\$)+ C " B"
 G ' .+ + + ('(I(%%
 L 8 *\$\$\$+ " \$% &
 *8 , 4 B@ 8 , +
 J : < *\$\$\$+ * 1 7 0
 *2 >R 2 FC+
 J 3T *\$(+ B 8
 + \$I (\$' - < (\$ /
 2 1 M E8
 J 3T C> @ *\$+ 1 - "
 *@ @ ?! / K >+
 J 92 *\$\$\$+ " ! "
 @ *17CF @ +
 J 92 *\$\$\$+ " ! "

< \$% , ,B 2 *E /
 K >+
 J L *%(+ (+ L J
 , *9 : 9 2 7V (I &&
 FR \$((+
 J C< 4, , 8 M *\$\$+) ' *
 *17CF @ +
 JB T7 *\$\$%+ -) ' * *F /
 K >+
 3, LE *\$\$+ # \$%
 *L M BN1 8 +
 F 8 *)'+ 9- + 7 , > * ,B 5 M
 5 < - , +
 F , 2 *\$)+ 7) \$ *M 9 "
 1 +
 F , 2 *\$&+ *M 9 /
 K >+
 F LM *\$\$)+ * \$% ! *@
 4 B @ @ /L+
 F < *\$%+ - *4 B
 8 @ 8 +
 /P @ *\$)+ # * " *C F 1 B
 @ <7+
 ? LF < , M 8 * + ' (\$) *17CF
 & (! < \$)
 @ +
 @ T *\$\$)+ 1 B G B
 ' 1(+)!
 @ 2 *\$ + * (" *@ 19 - "
 8 H /L+
 @ > C * + M ! B .+ -+ +
 1%%
 @ FLE *\$\$%+ - *8 , 4 "
 B @ 8 , +
 @ B LE *\$\$&+ 20 (" *8 "
 @ ?! +
 < , : @ *\$)+) *F 19
 / K >+
 < C < , : @ *\$)+)
 1 - *F 19 / K >+
 < 8 *\$%+ 2) *1 / K >+
 < ,, 9 *\$)((+ C > 1 0 + +
 %! \$
 < , M *\$)((+ T 7 ! + +
 + & !&%

< M *\$)%+ " 5 -
 *F 19 / K >+
 1 > 33 *\$ + 7 *L M BN 1
 / K >+
 1 J 8 B L 3 CM"8 *\$\$\$+) * +
) 7))
 *?! 4 B@ ?! +
 1 1 *\$)%+ / 0
 ** \$ 7 \$ ")
 + < - J 8 F *1 /
 K >+
 1 L *)&+ % 7 (* 1 3 +
 :
 1 CF 9 C < *\$\$\$+ \$
 *8 , 4 B@ 8 , +
 1 L8 * &(+ FP P P G P G
 5 5 5 ((5 * 3 4 5
 + + +)&I&
 1 L8 * &)%+ 1 PG HP P
 .+ + " + %I %
 1: W *\$(\$+ # * " *C F 1 B
 @ <7+
 5 39 *\$\$\$+ " 0 \$% & ()
 M 1 8 3 , < *8 , 4 B /
 K >+
 5 CF *\$'+ < "H ! 6 +. + +
 + +)I&
 5 3/ 2 E 777 *\$\$)+ 0 * *17CF
 @ +
 T E E *\$\$\$+ ' + 0 0+ + +
 7 \$ * 1 (*8 @ ?! +
 T < 1 *\$)%+ - ' (* @ 19 -
 8 H /L+
 M > L9 *\$)%+ - B !
 .+ +) + + I&&
 M > L9 *\$ + * * (" *8 @
 ?! 4 B@ / K >+
 K J *\$)+ 5 - *1 2 +
 K 5 *\$\$(+ 9 / I<
 ' 1(+ (&I((

Index

B D+ * *
#\$ >@ ??
\$ >@ ?? > > >>
8 >?
\$ >@ ?> >
8
*
@@
6
5 6 #
>
?
% 6, 1, @>
% 5 E
E , #
% || G 6 #>
% 6 , 5 , #?
, 5 , #* #
B D #
," / J 9
* @
8 @.
C9 " #&* #&

6 &
|| G 6 ??
8
\$ >?
\$.
8 #?
1 , *..
6 C , 6 >#
, 6 5 *?#
, J , 8
?
5 5 6 , *?&
5 *&*
*?> *?.
8 *@#
E *?
, 8 , *?> *?@ *&
1\$ 7 *&
" 1 *? *?. *&#
" *@
*@
|| 9 8
>
< = >
< = @*
?
? # >
?
- || : *@&
, 6 G , # >
, 6C , + G 6 >@ >
, J , *?
, , 6 #*? #*&
6 5 6 # ?*
"61 , @#
?*
>
1 ,
1 , *?? *&
>
1
@
\$, >. &
\$, 8 &
? .
**&
\$ # ? * #
1 , # ?
?

5 , ' & ## 9 8
 6 , , #? E # . ?
 6 # , #* # .
 1 # * #
 1 \$ *
 G , #? ##@
 #*
 #*
 # #
 , @.
 / G #
 G 5 , *
 (- # *
 B D # *
 / , @? ##& *?&
 / *&#
 / J G
 6 5 *?# *.>
 5 *#
 1 , *#*
 5 *** *#*
 , #*? #*#&
 E , ??
 K , # *#& # @
 M #?
 6 G #&
 8 #>?
 5 , #*?
 E ??
 5 #??
 K , #
 #*?
 6 *@@
 #*#* * *
 *#.
 #.#.
 C' , 1 ##

9 *.>
 6 *#&
 5 @
 F 1 , *@
 " *@
 F " , *? *?
 *@@
 6 C + , *? *?
 ? *@
 + 5 *@

 # #
 # #
 9 # #
 F " 8 *@*
 F " 1 , *@*
 F " E *@ *@* *@
 F " *@ *@*
 F " 6 *#& >
 F " * . @
 F G &&
 , 5 .>
 , 5 , *
 G &@
 G &@
 F 5 >
 F , # >
 F , 6 1
 # @
 F # >
 F , 5 , @ # *
 # # *
 F 6 5 *#? *?
 #.#.
 5 *#&
 F C , ?#
 O +1B D @?
 O 1 , @& *@
 O , , *
 O #.# &&
 #@
 O 8 & >@
 OI # G 6 ?#
 O 8 #>
 O R , #>#
 O # #>>

) , **
 \$ *
 C *>#
) 5 *
) E 6 >@ ?> >
 8 .
) 5 *# * @
 \$ *#&
) , > * .
 , , > *

6 >*
 1 , >>
) , (#
)
 , #
 , + #
 C' , 1 ##
 F G &&
 : .& .?
 C G #
 G &@
 G
 , .&
 8 #?
 #&
 || *
 +
) , >
) , 8 >*
) , #>&
) (#@
) #&@ @
 6 5
 O , #.& @
 : #. @
 @*
) #.
) 5
 1 8 &
 1 8
) 5 #??
)
 K , || # @
 , #
 , 5 , #
 5 , #*&
 5 , #
 # *
 K , 8 ##* * ?
 , " ?&
 : ^2 B D >
 : #. #
 : #. *
 : 6 @
 : G & 1
 : G >?
 : 6 ?*
 : , 5 , #
 : * @
 %\$ *
 6 * .

6 **
 , , 6 **&
 , , **
 8 , **
 , **
 1 , , 6 **.
 1 6 * .
 ** **>
 || **
 + , \$ 6 **#
 : , + , & ## # @ *#.
 : , + , # @
 A e > *#>
 : G .&
 : , G *
 : 8 ?
 5 &
 : 1 ,
 8 >
 1 8 .
 5 >*
 ' \$ 8 # #
 8 *@?
 8 @.
 #?
 , > . &
 6 @?
 O &
 ?
 ' \$ 8 # #
 @@
 #*.
 >*
 5 E .& .. @&
 8 >
 6 , .& 6 @?
 6 , .&
 8
 ?
 , &
 8 1 ,
 " 6 @
 : M &? & . #? *
 : 8 > .
 : \$??
 : \$??
 : \$?>
 9 # #
 ?
 : 8) E 6
 : 8 *?> *?@ * &
 , 1 , *?? * &

' (. # ## ? ##* H 8 #*.
H 6 >@ ? &&
' .? , , ?
' 8 8 * + ?@ &@
' " " || G 6 ? H 1
' E * . 5 #&
' , 6 * . 5 , 1 1 G #
' 8 @@ #*.
' # H 6 ?>
' C 5 & H , * ** *
. # , , **
\$ 8 6 @. #.
> & , > *
N # > @ @
5 # @ # ##? 5 > @ @>
, * ##? #* #?*
5 5 *# # * , -C9 G 6 #
G ##. 5 E 8 .& @&
C G # .. 8 >
, 5 , . , 8 >
N * M #?
+ 1 #?
\$ >@ ?? M 1 , &? & #?*
\$ >@ ?? > > >> M , , 5 , #? ##@
\$ >@ ?> > M)
9 6 *@@ sym .&
, # # ?
, > , > .
\$? G .&
5 5 *.& 6 G .&
5 , 8 > . 6 G #&
G &? , 8 C + , *..
H \$ *#& , 5 , ##& ##? #?
, * # + 8 * @*
5 , * *#& * * + * @ * @*
1 *#& * N #. # @ # @ #&
, * ?
5 * .
E * .
1 , , 6 * * C * * * * >
C * * * * >
+ *#& , , 1 * .
H ?* , 6 , B/) D * > *
H # , * @ 5 * .
H , @ 5 * .
H 1 , , 6 E * .
\$ * * , 5 >
H , , 5 , ?
H 5 , #*? ' > 1\$ 7 *&

1 + @>
 1 M #?>
 , 5 , ##
 5 , #>
 || * @
 #
 >
 # ?
 , 5 , ## ##*
 G #
 || ##.
 8 ##?
 5
 , E ?&
 5 *.&
 5 , @? *.&
 1 G " >? >>
 G &
 5 6 &#
 @ *@
 , , @.
 O , , * @.
 * #
 , , @.
 " @*
 1 , @&
 @*
 6 @?
 , , @.
 @@
 , , @.
 6 *
 6 6 * &
 E8 #
 J H/ 6 * >
 J 8 *@?
 , 6 6 8 ! || G 6 ?#
 @?
 G , #>.
 C: 5 *&*

5 *@
 , 5 , #*
 8 ?
 , 5 , #*
 *@
 6 1 , *..
 6 , 8 .&
 6 ||
 5 5
 @> *?*
 , #*?
 +
 @
 8 @* *?* *?& *&#
 & *@>
 5 #>?
 1 , @
 #
 1 6 , 8
 #>
 || *#.
 **&
 \$ *#&
 8 ?
 8 &
 (*.>
 (, >.
 # \$ >@
 \$ >@
 \$ >@
 \$?
 " 1 *@
 " *@
 1 , >> ? & . *
 &&
 8 ! || G 6 ?#
 Y \$ 6 **#
 **>