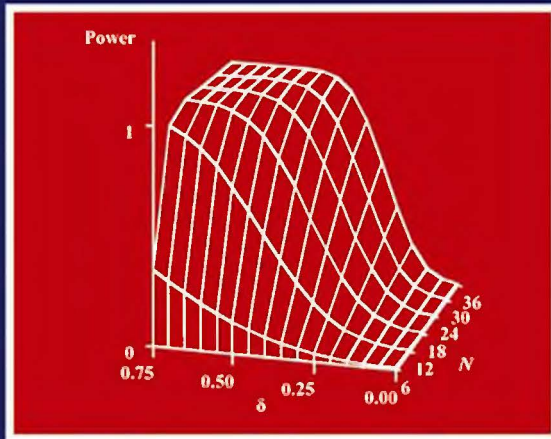


Linear Model Theory

Univariate, Multivariate, and Mixed Models



Keith E. Muller
Paul W. Stewart



LINEAR MODEL THEORY

Univariate, Multivariate, and Mixed Models

Keith E. Muller

*University of North Carolina
Department of Biostatistics
Chapel Hill, NC*

Paul W. Stewart

*University of North Carolina
Department of Biostatistics
Chapel Hill, NC*



A JOHN WILEY & SONS, INC., PUBLICATION

LINEAR MODEL THEORY

LINEAR MODEL THEORY

Univariate, Multivariate, and Mixed Models

Keith E. Muller

*University of North Carolina
Department of Biostatistics
Chapel Hill, NC*

Paul W. Stewart

*University of North Carolina
Department of Biostatistics
Chapel Hill, NC*

 **WILEY-
INTERSCIENCE**

A JOHN WILEY & SONS, INC., PUBLICATION

Copyright © 2006 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.
Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permission>.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic format. For information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Muller, Keith E.

Linear model theory : univariate, multivariate, and mixed models / Keith E. Muller,
Paul W. Stewart.
p. cm.

Includes bibliographical references and index.

ISBN-10 0-471-21488-4

ISBN-13 978-0-471-21488-5

I. Linear models (Statistics). I. Stewart, Paul Wilder. II. Title.

QA279.M86 2006

519.5—dc22

2006044266

Printed in the United States of America.

10 9 8 7 6 5 4 3 2 1

Contents

Preface

xiii

1 Matrix Algebra for Linear Models

1

- 1.1 Notation
- 1.2 Some Operators and Special Types of Matrices
- 1.3 Five Kinds of Multiplication
- 1.4 The Direct Sum
- 1.5 Rules of Operations
- 1.6 Other Special Types of Matrices
- 1.7 Quadratic and Bilinear Forms
- 1.8 Vector Spaces and Rank
- 1.9 Finding Rank
- 1.10 Determinants
- 1.11 The Inverse and Generalized Inverse
- 1.12 Eigenanalysis (Spectral Decomposition)
- 1.13 Some Factors of Symmetric Matrices
- 1.14 Singular Value Decomposition
- 1.15 Projections and Other Functions of a Design Matrix
- 1.16 Special Properties of Patterned Matrices
- 1.17 Function Optimization and Matrix Derivatives
- 1.18 Statistical Notation Involving Matrices
- 1.19 Statistical Formulas
- 1.20 Principal Components
- 1.21 Special Covariance Patterns

2 The General Linear Univariate Model	39
2.1 Motivation	
2.2 Model Concepts	
2.3 The General Linear Univariate Linear Model	
2.4 The Univariate General Linear Hypothesis	
2.5 Tests about Variances	
2.6 The Role of the Intercept	
2.7 Population Correlation and Strength of Relationship	
2.8 Statistical Estimates	
2.9 Testing the General Linear Hypothesis	
2.10 Confidence Regions for θ	
2.11 Sufficient Statistics for the Univariate Model	
Exercises	
3 The General Linear Multivariate Model	55
3.1 Motivation	
3.2 Definition of the Multivariate Model	
3.3 The Multivariate General Linear Hypothesis	
3.4 Tests About Covariance Matrices	
3.5 Population Correlation	
3.6 Statistical Estimates	
3.7 Overview of Testing Multivariate Hypotheses	
3.8 Computing MULTIREP Tests	
3.9 Computing UNIREP Tests	
3.10 Confidence Regions for Θ	
3.11 Sufficient Statistics for the Multivariate Model	
3.12 Allowing Missing Data in the Multivariate Model	
Exercises	
4 Generalizations of the Multivariate Linear Model	79
4.1 Motivation	
4.2 The Generalized General Linear Univariate Model: Exact and Approximate Weighted Least Squares	
4.3 Doubly Multivariate Models	
4.4 Seemingly Unrelated Regressions	
4.5 Growth Curve Models (GMANOVA)	
4.6 The Relationship of the GCM to the Multivariate Model	
4.7 Mixed, Hierarchical, and Related Models	

5 The Linear Mixed Model	91
5.1 Motivation	
5.2 Definition of the Mixed Model	
5.3 Distribution-Free and Noniterative Estimates	
5.4 Gaussian Likelihood and Iterative Estimates	
5.5 Tests about β (Means, Fixed Effects)	
5.6 Tests of Covariance Parameters, τ (Random Effects)	
Exercises	
6 Choosing the Form of a Linear Model for Analysis	101
6.1 The Importance of Understanding Dependence	
6.2 How Many Variables per Independent Sampling Unit?	
6.3 What Types of Variables Play a Role?	
6.4 What Repeated Sampling Scheme Was Used?	
6.5 Analysis Strategies for Multivariate Data	
6.6 Cautions and Recommendations	
6.7 Review of Linear Model Notation	
7 General Theory of Multivariate Distributions	115
7.1 Motivation	
7.2 Notation and Concepts	
7.3 Families of Distributions	
7.4 Cumulative Distribution Function	
7.5 Probability Density Function	
7.6 Formulas for Probabilities and Moments	
7.7 Characteristic Function	
7.8 Moment Generating Function	
7.9 Cumulant Generating Function	
7.10 Transforming Random Variables	
7.11 Marginal Distributions	
7.12 Independence of Random Vectors	
7.13 Conditional Distributions	
7.14 (Joint) Moments of Multivariate Distributions	
7.15 Conditional Moments of Distributions	
7.16 Special Considerations for Random Matrices	

8	Scalar, Vector, and Matrix Gaussian Distributions	139
8.1	Motivation	
8.2	The Scalar Gaussian Distribution	
8.3	The Vector (“Multivariate”) Gaussian Distribution	
8.4	Marginal Distributions	
8.5	Independence	
8.6	Conditional Distributions	
8.7	Asymptotic Properties	
8.8	The Matrix Gaussian Distribution	
8.9	Assessing Multivariate Gaussian Distribution	
8.10	Tests for Gaussian Distribution	
	Exercises	
9	Univariate Quadratic Forms	169
9.1	Motivation	
9.2	Chi-Square Distributions	
9.3	General Properties of Quadratic Forms	
9.4	Properties of Quadratic Forms in Gaussian Vectors	
9.5	Independence among Linear and Quadratic Forms	
9.6	The ANOVA Theorem	
9.7	Ratios Involving Quadratic Forms	
	Exercises	
10	Multivariate Quadratic Forms	193
10.1	The Wishart Distribution	
10.2	The Characteristic Function of the Wishart	
10.3	Properties of the Wishart	
10.4	The Inverse Wishart	
10.5	Related Distributions	
	Exercises	
11	Estimation for Univariate and Weighted Linear Models	209
11.1	Motivation	
11.2	Statement of the Problem	
11.3	(Unrestricted) Linearly Equivalent Linear Models	
11.4	Estimability and Criteria for Checking It	
11.5	Coding Schemes and the Essence Matrix	
11.6	Unrestricted Maximum Likelihood Estimation of β	
11.7	Unrestricted BLUE Estimation of β	
11.8	Unrestricted Least Squares Estimation of β	

- 11.9 Unrestricted Maximum Likelihood Estimation of θ
- 11.10 Unrestricted BLUE Estimation of θ
- 11.11 Related Distributions
- 11.12 Formulations of Explicit Restrictions of β and θ
- 11.13 Restricted Estimation Via Equivalent Models
- 11.14 Fitting Piecewise Polynomial Models Via Splines
- 11.15 Estimation for the GGLM: Weighted Least Squares
Exercises

12 Estimation for Multivariate Linear Models **243**

- 12.1 Alternate Formulations of the Model
- 12.2 Estimability in the Multivariate GLM
- 12.3 Unrestricted Likelihood Estimation
- 12.4 Estimation of Secondary Parameters
- 12.5 Estimation with Multivariate Restrictions
- 12.6 Unrestricted Estimation With Compound Symmetry:
the “Univariate” Approach to Repeated Measures
Exercises

13 Estimation for Generalizations of Multivariate Models **263**

- 13.1 Motivation
- 13.2 Criteria and Algorithms
- 13.3 Weighted Estimation of B and Σ
- 13.4 Transformations among Growth Curve Designs
- 13.5 Within-Individual Design Matrices
- 13.6 Estimation Methods
- 13.7 Relationships to the Univariate and Mixed Models
Exercises

14 Estimation for Linear Mixed Models **279**

- 14.1 Motivation
- 14.2 Statement of the General Linear Mixed Model
- 14.3 Estimation and Estimability
- 14.4 Some Special Types of Models
- 14.5 ML Estimation
- 14.6 REML Estimation
- 14.7 Small-Sample Properties of Estimators
- 14.8 Large-Sample Properties of Variance Estimators
- 14.9 Conditional Estimation of d_i and BLUP Prediction
Exercises

15 Tests for Univariate Linear Models	289
15.1 Motivation	
15.2 Testability of Univariate Hypotheses	
15.3 Tests of a Priori Hypotheses	
15.4 Related Distributions	
15.5 Transformations and Invariance Properties	
15.6 Confidence Regions for θ	
Exercises	
16 Tests for Multivariate Linear Models	311
16.1 Motivation	
16.2 Testability of Multivariate Hypotheses	
16.3 Tests of a Priori Hypotheses	
16.4 Linear Invariance	
16.5 Four Multivariate Test Statistics	
16.6 Which Multivariate Test Is Best?	
16.7 Univariate Approach to Repeated Measures: UNIREP	
16.8 More on Invariance Properties	
16.9 Tests of Hypotheses about Σ	
16.10 Confidence Regions for Θ	
Exercises	
17 Tests for Generalizations of Multivariate Linear Models	337
17.1 Motivation	
17.2 Doubly Multivariate Models	
17.3 Missing Responses in Multivariate Linear Models	
17.4 Exact and Approximate Weighted Least Squares	
17.5 Seemingly Unrelated Regressions	
17.6 Growth Curve Models (GMANOVA)	
17.7 Testing Hypotheses in the GCM	
17.8 Confidence Bands for Growth Curves	
18 Tests for Linear Mixed Models	341
18.1 Overview	
18.2 Estimability of $\theta = C\beta$	
18.3 Likelihood Ratio Tests of $C\beta$	
18.4 Likelihood Ratio Tests Involving τ	
18.5 Test Size of Wald-Type Tests of β Using REML	
18.6 Using Wald-Type Tests of β with REML	
18.7 Using Wald-Type Tests of $\{\beta, \tau\}$ with REML	

19 A Review of Multivariate and Univariate Linear Models	349
19.1 Matrix Gaussian and Wishart Properties	
19.2 Design Matrix Properties	
19.3 Model Components	
19.4 Primary Parameter and Related Estimators	
19.5 Secondary Parameter Estimators	
19.6 Added-Last and Added-in-Order Tests	
20 Sample Size for Univariate Linear Models	361
20.1 Sample Size Consulting: Before You Begin	
20.2 The Machinery of a Power Analysis	
20.3 Independent t Example	
20.4 Paired t Example	
20.5 The Impact of Using $\hat{\sigma}^2$ or $\hat{\beta}$ in Power Analysis	
20.6 Random Predictors	
20.7 Internal Pilot Designs	
20.8 Other Criteria for Choosing a Sample Size	
Exercises	
21 Sample Size for Multivariate Linear Models	371
21.1 The Machinery of a Power Analysis	
21.2 Paired t Example	
21.3 Time by Treatment Example	
21.4 Comparing between and within Designs	
21.5 Some Invariance Properties	
21.6 Random Predictors	
21.7 Internal Pilot Designs	
Exercises	
22 Sample Size for Generalizations of Multivariate Models	383
22.1 Motivation	
22.2 Sample Size Methods for Growth Curve Models	
23 Sample Size for Linear Mixed Models	385
23.1 Motivation	
23.2 Methods	
23.3 Internal Pilot Designs	

Appendix: Computing Resources

387

References

393

Index

405

Preface

MOTIVATION FOR THE BOOK

Statisticians often use linear models for data analysis and for developing new statistical methods. Success in either endeavor requires a solid understanding of the underlying theory. Historically, univariate, multivariate, and mixed linear models have been discussed separately. In contrast, we give a *unified* treatment in order to make clear the distinctions among the three classes of models. No single model class proves uniformly best. Therefore choosing the best approach requires a detailed knowledge of the differences and similarities.

A student needs to acquire four sets of skills. (1) Using all three classes of linear models correctly in practice requires knowing enough theory, especially a deep understanding of assumptions and their possible violations. However, we leave detailed discussion of diagnostics to others. (2) Correct use also requires knowing when to choose one type over another. (3) Understanding the theory helps guide when *not* to use any of the models. (4) Finally, developing new methods requires a detailed knowledge of known work.

TOPICS COVERED

We focus on linear models of interval scale responses with finite second moments, especially models with correlated observations and Gaussian errors. Such correlations always create additional complexity. In contrast to most “multivariate” books, most classical techniques, including cluster analysis, factor analysis, discriminant analysis, and canonical correlation, receive little attention.

Meeting our goals in a book appropriate for a one-semester course required omitting many worthwhile topics. Even so, the book includes more material than usually covered entirely in a four-credit class. On the other hand, it may seem wasteful to include separate and overlapping treatments of univariate, multivariate, and mixed models. However, our students have adamantly preferred the current organization. We accept some duplication for the sake of clarity.

NOTATION

We sought a precise but accessible presentation of the theory underlying practice, illustrated with examples. Fairness to our students and readers required

creating compatible univariate, multivariate, and mixed model notation. Many authors describe a univariate model in terms of a single observation. Others describe univariate and multivariate models in matrix notation for all observations and all independent sampling units. In contrast, mixed model discussions typically describe the observations for only one independent sampling unit. Table 6.5 serves as a Rosetta stone by allowing translations between models and publications. Discussing all three forms of all models clarifies differences and helps avoid many sources of confusion.

PREREQUISITES

We developed the book for a four-hour class required of our doctoral students. Most have the equivalent of a Master's Degree in Biostatistics, including two semesters of probability and inference, at the level of Wackerly, Mendenhall, and Scheaffer (1996). Students also need a solid background in applied univariate linear models from a matrix perspective, as in Muller and Fetterman (2002).

Most explanations and proofs use matrix algebra. Rank, basis space, eigenanalysis, the singular value decomposition, and generalized inverses play central roles. Although we summarize key matrix results in Chapter 1, the reader without sufficient background should invest time in studying one of the many good matrix theory books, such as Schott (2005). Section 1.15 merits special and repeated attention in understanding and proving many results in estimation and inference. Prerequisite material in sections 1.11-1.14 also may require study.

We use the most basic properties of complex arithmetic to allow simplifying some proofs by using characteristic functions. Measure theoretic and contour integration methods are mentioned only occasionally, and not required.

ACKNOWLEDGMENTS

We owe a great debt to the authors of many earlier books about linear models. We give explicit citations for particular results and also for conceptual approaches. We wrote the book while the first author was a Professor in the Department of Biostatistics at the University of North Carolina Hill.

Many colleagues helped us along the way. We especially thank Ronald W. Helms, Lloyd J. Edwards, and Christopher S. Coffey for their guidance and friendship. Teaching assistants Stacey Major, William K. Pan, Hae-Young Kim, and J. (Chris) Slaughter, as well as a number of graders, helped shape the organization and choice of topics by providing a student's perspective. Our hard-working and enthusiastic students motivated us to start the book and kept us going to the finish. As in all such endeavors, the book would not exist without the love, support, and forbearance of our families, especially our wives, Sally and Dawn.

CHAPTER 1

Matrix Algebra for Linear Models

1.1 NOTATION

Graybill (1969), Searle (1982), Harville (1996), and Schott (2005) provided thorough introductions to matrix algebra for statistics. We summarize only key results here. Substantial omissions include the deletion of nearly all proofs, as well as consideration of more general forms not commonly used in statistics. Also, many issues of numerical accuracy have been ignored. Some of the formulas described here, although very useful for understanding concepts, prove numerically unstable with typical computer precision.

Braces, $\{ \}$, indicate sets and *brackets*, $[]$, indicate matrices or vectors (arrays). In a distinct use of the same symbols, mathematical expressions will be grouped by using the nesting sequence $\{[()]\}$, which may be iterated as $\{[(\{[()]\})]\}$.

Definition 1.1 (a) A *matrix* is a rectangular, two-dimensional array of *elements*. Writing $\mathbf{A} = \{a_{ij}\}$ says \mathbf{A} is the matrix with element a_{ij} at row i and column j . Here i is the row index, while j is the column index, which are always written in row-column order.

(b) A *vector* is any matrix with exactly one column, such as $\mathbf{v} = \begin{bmatrix} 7 \\ 8 \end{bmatrix}$.

(c) A *scalar*, such as $s = 6$, can be expressed as a vector with one row or as a matrix with one row and one column, written $s = \mathbf{s} = \mathbf{S}$.

We restrict attention to real numbers and finite dimensions. With r rows, c columns, \mathbf{A} is $r \times c$ (r by c):

$$\mathbf{A} = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1c} \\ a_{21} & & & \vdots \\ \vdots & & & \vdots \\ a_{r1} & \cdots & \cdots & a_{rc} \end{bmatrix}. \quad (1.1)$$

In turn, $\langle \mathbf{A} \rangle_{ij} = a_{ij}$ indicates element i, j has been extracted from \mathbf{A} .

Although not all authors do so, we are scrupulous about the distinction between a matrix of one row, such as $\mathbf{A} = [a_1 \ a_2]$, and a vector, $\mathbf{b} = \begin{bmatrix} a_1 \\ a_1 \end{bmatrix}$. The vector \mathbf{b}

can also be written in terms of a transpose (defined in the next section), $\mathbf{b} = \mathbf{A}'$. Doing so not only avoids notational ambiguity, but also builds in many consistency checks by requiring dimensions and symbol types (\mathbf{a} or \mathbf{A}) to align.

As an aid to working with matrices, we always use bold typeface in word processing software, as in the present book. Most, but not all, statistical journals require the convention. With the convention, a represents a scalar, \mathbf{a} represents a vector, and \mathbf{A} represents a matrix. When handwriting expressions, we highly recommend always putting a tilde or dash under any matrix or vector to indicate boldface. It is *extremely* helpful to write the dimensions of each matrix in an equation underneath the equation. Using the transpose, matrix multiplication, and inverse operators (introduced later in the chapter) illustrates the idea, with

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (1.2)$$

being much less informative than

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \quad (1.3)$$

$p \times 1 \quad [(p \times n)(n \times p)](p \times n)(n \times 1)$

The practice saves a great deal of time (otherwise spent being confused). More bluntly, if one does not know the dimensions, one cannot understand the equation.

We reserve superscripts for operators and use subscripts for descriptors, such as in x , x^2 , x_2 . Often, functional notation, such as $x(c, \alpha)$, provides a better alternative than a long and elaborate subscript descriptor.

1.2 SOME OPERATORS AND SPECIAL TYPES OF MATRICES

Definition 1.2 (a) A *square* matrix has the same number of rows as columns.

(b) For \mathbf{A} ($r \times c$), the (main) *diagonal* of \mathbf{A} is $\{a_{11}, a_{22}, \dots, a_{cc}\}$.

(c) A square matrix is *diagonal* if all elements off the main diagonal are zero; if $i \neq j$, then $a_{ij} = 0$, while a_{ii} can be anything.

Definition 1.3 Writing $\text{Dg}(\mathbf{v}) = \text{Dg}(\{v_j\})$ indicates creating a square diagonal matrix from a vector or elements of a set, as in

$$\text{Dg} \left(\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \right) = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & 0 \\ 0 & 0 & 3 \end{bmatrix}. \quad (1.4)$$

Definition 1.4 A prime indicates *transpose*. If $\mathbf{A} = \{a_{ij}\}$ is $r \times c$, then $\mathbf{A}' = \{a_{ji}\}$ is $c \times r$. The transpose operation causes rows to become columns and columns to become rows.

Definition 1.5 A *symmetric* matrix is a square matrix such that $a_{ij} = a_{ji}$. Equivalently, $\mathbf{A}' = \mathbf{A}$.

As mentioned earlier, we use the term “vector” only for $n \times 1$ arrays and never for $1 \times n$ arrays. A $1 \times n$ array will always be written as a matrix, such as $\mathbf{A} = [1 \ 2 \ 3]$, or as a transposed vector, $\mathbf{b}' = [1 \ 2 \ 3]$. Here $\mathbf{A} = \mathbf{b}'$ and $\mathbf{A}' = \mathbf{b}$.

Definition 1.6 An *identity* matrix, \mathbf{I} or \mathbf{I}_n , is a square matrix with all 1's on the main diagonal, and all 0's off-diagonal. Equivalently, $a_{ij} = 0$ if $i \neq j$ and $a_{ij} = 1$ if $i = j$.

Definition 1.7 A *zero* matrix, $\mathbf{0}$, has $a_{ij} \equiv 0$ and may be written $\mathbf{0}_r$ to indicate a vector or $\mathbf{0}_{r \times c}$ to indicate a matrix for clarity.

Similarly an $n \times 1$ vector with all elements 1 is written $\mathbf{1}$ or $\mathbf{1}_n = [1 \ 1 \ \dots \ 1]'$. Also, $\mathbf{1}_n \mathbf{1}'_n$ is an $n \times n$ matrix of all 1's.

Definition 1.8 An *upper triangular* matrix has $a_{ij} = 0$ for $i > j$, such as

$$\mathbf{U} = \begin{bmatrix} 1 & 5 & 6 \\ 0 & 2 & 4 \\ 0 & 0 & 3 \end{bmatrix}. \tag{1.5}$$

A *lower triangular* matrix has $a_{ij} = 0$ for $i < j$:

$$\mathbf{L} = \begin{bmatrix} 1 & 0 & 0 \\ 7 & 2 & 0 \\ 8 & 9 & 3 \end{bmatrix}. \tag{1.6}$$

Definition 1.9 A *partitioned* matrix (supermatrix) has elements grouped meaningfully by combinations of vertical and horizontal slicing, indicated $\mathbf{A} = \{\mathbf{A}_{jk}\}$. Necessarily \mathbf{A}_{jk} and $\mathbf{A}_{j'k'}$ have the same number of rows, while \mathbf{A}_{jk} and $\mathbf{A}_{j'k}$ have the same number of columns.

Definition 1.10 A *block diagonal* matrix is a partitioned matrix with all partitions zero except possibly $\{\mathbf{A}_{jj}\}$.

Two examples are the block diagonal matrix

$$\mathbf{A} = \begin{bmatrix} 1 & 2 & \vdots & 0 & 0 \\ 3 & 4 & \vdots & 0 & 0 \\ \dots & \dots & \ddots & \dots & \dots \\ 0 & 0 & \vdots & 5 & 6 \\ 0 & 0 & \vdots & 7 & 8 \end{bmatrix} \tag{1.7}$$

and the general partitioned matrix

$$B = \begin{bmatrix} 1 & \vdots & 0 & 0 \\ 1 & \vdots & 0 & 0 \\ 2 & \vdots & 0 & 0 \\ 3 & \vdots & 0 & 0 \end{bmatrix}. \quad (1.8)$$

Conformation of shapes requires consistent partitioning. If

$$A = \begin{bmatrix} B & C & D \\ E & F & G \end{bmatrix} \quad (1.9)$$

then B , C , and D have the same number of rows, while B and E have the same number of columns, etc. However, complete uniformity of dimensions is not required (the number of rows of B need not equal the number of rows of E). It would be hard to overemphasize the value of partitioned matrices in deriving algebraic and statistical properties for linear models. Expressions can often be greatly simplified by taking advantage of special properties of partitioned matrices for basic operations (matrix summation, multiplication, etc.) and more complicated operations (determinants, inverses, etc.).

Definition 1.11 For $r \times c$ A , writing $\text{col}_k(A) = \mathbf{a}_k$ indicates extracting $r \times 1$ column k from A . Writing $A_j = \text{row}_j(A)$ indicates extracting a particular $1 \times c$ row.

As an important example of partitioning, $r \times c$ A can be expressed in terms of its c column vectors, $\{\mathbf{a}_j\}$, with \mathbf{a}_j of dimension $r \times 1$, or its r rows, $\{A_k\}$, with A_k of dimension $c \times 1$. In summary,

$$\begin{aligned} A &= [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_c] \\ &= \begin{bmatrix} A_1 \\ A_2 \\ \vdots \\ A_r \end{bmatrix}. \end{aligned} \quad (1.10)$$

Definition 1.12 (a) Writing $A = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_c]$, which requires $\{\mathbf{a}_j\}$ to be $r \times 1$, indicates $\{\mathbf{a}_j\}$ have been *horizontally concatenated*.

(b) Writing $A = \begin{bmatrix} A_1 \\ \vdots \\ A_r \end{bmatrix}$, which requires $\{A_k\}$ to be $c \times 1$, indicates $\{A_k\}$ have been *vertically concatenated*.

Definition 1.13 Writing $\text{vec}(\)$ indicates all elements of a matrix have been stacked *by column*, as in $\mathbf{b}_1 = \text{vec}(A)$, because it creates an $(rc) \times 1$ vector from an $r \times c$ matrix. Equivalently, the columns have been vertically concatenated.

If $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \cdots \ \mathbf{a}_c]$ then

$$\text{vec}(\mathbf{A}) = \begin{bmatrix} \mathbf{a}_1 \\ \mathbf{a}_2 \\ \vdots \\ \mathbf{a}_c \end{bmatrix}. \tag{1.11}$$

In turn, $\mathbf{b}_2 = \text{vec}(\mathbf{A}')$ also creates an $(rc) \times 1$ vector which differs from \mathbf{b}_1 only by permutation of the rows. Creating $\text{vec}(\mathbf{A}')$ stacks the matrix by rows:

$$\text{vec}(\mathbf{A}') = \begin{bmatrix} \mathbf{A}'_1 \\ \mathbf{A}'_2 \\ \vdots \\ \mathbf{A}'_r \end{bmatrix}. \tag{1.12}$$

If $r = c$ and $\mathbf{A} = \mathbf{A}'$ then only $r(r + 1)/2$ elements are distinct. The r^2 elements are not functionally independent. The $\text{vech}()$ operator extracts the distinct elements into a vector:

$$\mathbf{z} = \text{vech} \left(\begin{bmatrix} a & b & c \\ b & d & e \\ c & e & f \end{bmatrix} \right) = \begin{bmatrix} a \\ b \\ c \\ d \\ e \\ f \end{bmatrix}. \tag{1.13}$$

Definition 1.14 The *trace* of an $n \times n$ (square) matrix is $\text{tr}(\mathbf{A}) = \sum_{i=1}^n a_{ii}$.

Definition 1.15 Matrices *conform* for an operation if their sizes allow the result of the operation to exist. Matrices do not conform for an operation if their dimensions do not allow the desired operation.

Definition 1.16 *Matrix addition* yields $\mathbf{A} + \mathbf{B} = \{a_{ij} + b_{ij}\}$ while *matrix subtraction* yields $\mathbf{A} - \mathbf{B} = \{a_{ij} - b_{ij}\}$. Either result exists only if \mathbf{A} and \mathbf{B} are the same size (and thereby conform for the operation).

1.3 FIVE KINDS OF MULTIPLICATION

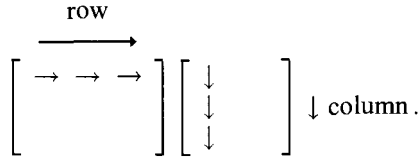
Definition 1.17 *Scalar multiplication* of a matrix gives $\mathbf{A}b = b\mathbf{A} = \{ba_{ij}\}$.

Definition 1.18 If \mathbf{A} and \mathbf{B} are both $r \times c$, then *elementwise multiplication* gives $\mathbf{A}\#\mathbf{B} = \{a_{ij}b_{ij}\} = \mathbf{C}$, with \mathbf{C} also $r \times c$.

Definition 1.19 (a) Matrix multiplication of \mathbf{A} ($r \times c$) and \mathbf{B} ($c \times d$) gives $\mathbf{AB} = \mathbf{C} = \{c_{jk}\}$ for $c_{jk} = \sum_{m=1}^c a_{km}b_{mk}$, with \mathbf{C} $r \times d$.
(b) If \mathbf{A} is $r \times r$, the matrix power of \mathbf{A} is $\mathbf{A}^k = \mathbf{A}_1\mathbf{A}_2 \cdots \mathbf{A}_k$, with $\mathbf{A}_j = \mathbf{A}$.

Definition 1.20 (a) The cross product of $r \times 1$ \mathbf{a} and \mathbf{b} is $\mathbf{a}'\mathbf{b} = \sum_{k=1}^r a_k b_k$.
(b) The dot product of $r \times 1$ \mathbf{a} and \mathbf{b} is $\mathbf{a} \cdot \mathbf{b} = \cos(\theta)\sqrt{\mathbf{a}'\mathbf{a}\mathbf{b}'\mathbf{b}}$, with θ the angle between the two vectors. Necessarily $\mathbf{a}'\mathbf{b} = 0$ if and only if $\theta = 90^\circ$.
(c) If \mathbf{a} is $r \times 1$, then $\mathbf{a}'\mathbf{a} = \sum_{k=1}^r a_k^2$ is the inner product of the vector.
(d) If \mathbf{a} is $r \times 1$, then $\mathbf{a}\mathbf{a}' = \{a_i a_j\}$ is the outer product of the vector.

Matrix multiplication can be expressed as a collection of cross products. Multiplying row j of \mathbf{A} with column j of \mathbf{B} yields $c_{jk} = \{\text{row}_j(\mathbf{A})\text{col}_k(\mathbf{B})\}$:



Lemma 1.1 (a) Premultiplying by a (square) diagonal matrix scales the rows, and postmultiplying by a (square) diagonal matrix scales the columns.
(b) The result generalizes to partitioned matrices with conforming partitions.

In particular, for conforming matrices $\mathbf{IA} \equiv \mathbf{A}$ and $\mathbf{AI} \equiv \mathbf{A}$. For a 2×3 matrix, \mathbf{A} , multiplication by diagonal matrices gives

$$\mathbf{GA} = \mathbf{GAI} = \begin{bmatrix} g_1 & 0 \\ 0 & g_2 \end{bmatrix} \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} g_1 a & g_1 b & g_1 c \\ g_2 d & g_2 e & g_2 f \end{bmatrix} \quad (1.14)$$

$$\mathbf{AH} = \mathbf{IAH} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} a & b & c \\ d & e & f \end{bmatrix} \begin{bmatrix} h_1 & 0 & 0 \\ 0 & h_2 & 0 \\ 0 & 0 & h_3 \end{bmatrix} = \begin{bmatrix} ah_1 & bh_2 & ch_3 \\ dh_1 & eh_2 & fh_3 \end{bmatrix} \quad (1.15)$$

Definition 1.21 The horizontal direct product creates a new matrix by elementwise multiplication of pairs of columns from two matrices with the

same number of rows:

$$\begin{aligned}
 \mathbf{A} \odot \mathbf{B} &= \begin{bmatrix} a & b \\ c & d \\ e & f \end{bmatrix} \odot \begin{bmatrix} r & s & t \\ u & v & w \\ x & y & z \end{bmatrix} \\
 &= \begin{bmatrix} ar & as & at & br & bs & bt \\ cu & cv & cw & du & dv & dw \\ ex & ey & ez & fx & fy & fz \end{bmatrix}. \tag{1.16}
 \end{aligned}$$

To operate on rows (1) transpose each operand, (2) compute the product, and (3) transpose the result, as with $(\mathbf{A}' \odot \mathbf{B}')'$. The operator could be called the vertical or column direct product.

Definition 1.22 The *direct* (or Kronecker) product is

$$\begin{aligned}
 \mathbf{A} \otimes \mathbf{B} &= \{a_{ij}\mathbf{B}\} \\
 &= \begin{bmatrix} a_{11}\mathbf{B} & \cdots & a_{1c}\mathbf{B} \\ a_{21}\mathbf{B} & \ddots & \vdots \\ \vdots & & \vdots \\ a_{r1}\mathbf{B} & \cdots & a_{rc}\mathbf{B} \end{bmatrix}. \tag{1.17}
 \end{aligned}$$

With $r \times c$ \mathbf{A} and $s \times d$ \mathbf{B} , the result has dimension $(rs) \times (cd) = (\text{rows} \times \text{columns})$. Some authors choose to define $\{\mathbf{A}b_{ij}\}$ as the direct product, which produces a different matrix.

1.4 THE DIRECT SUM

Definition 1.23 The *direct sum* operator creates a block diagonal matrix from any set of square matrices:

$$\begin{aligned}
 \bigoplus_{j=1}^J \mathbf{A}_j &= \mathbf{A}_1 \oplus \mathbf{A}_2 \oplus \cdots \oplus \mathbf{A}_J \\
 &= \begin{bmatrix} \mathbf{A}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{A}_2 & \mathbf{0} & \vdots \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{A}_J \end{bmatrix}. \tag{1.18}
 \end{aligned}$$

Lemma 1.2 In general $\text{tr}(\mathbf{A}_1 \oplus \mathbf{A}_2) = \text{tr}(\mathbf{A}_1) + \text{tr}(\mathbf{A}_2)$. If \mathbf{A}_j and \mathbf{B}_j are both $n_j \times n_j$, then

$$\left(\bigoplus_{j=1}^J \mathbf{A}_j \right) + \left(\bigoplus_{j=1}^J \mathbf{B}_j \right) = \bigoplus_{j=1}^J (\mathbf{A}_j + \mathbf{B}_j) \quad (1.19)$$

$$\left(\bigoplus_{j=1}^J \mathbf{A}_j \right) \left(\bigoplus_{j=1}^J \mathbf{B}_j \right) = \bigoplus_{j=1}^J (\mathbf{A}_j \mathbf{B}_j). \quad (1.20)$$

Any direct product of the form $\mathbf{I} \otimes \mathbf{B}$ is a special case of the direct sum:

$$\mathbf{I}_J \otimes \mathbf{B} = \bigoplus_{j=1}^J \mathbf{B}. \quad (1.21)$$

Direct products including an identity matrix, $\mathbf{A} \otimes \mathbf{I}$ or $\mathbf{I} \otimes \mathbf{B}$, occur often in expressions for covariance matrices of data in clusters of fixed size. A common form occurs in describing the covariance matrix of data observed in N clusters of constant size, with homogeneity of covariance between clusters:

$$\begin{aligned} \Xi &= \mathbf{I}_N \otimes \Sigma \\ &= \begin{bmatrix} \Sigma & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma & \mathbf{0} & \vdots \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \Sigma \end{bmatrix}. \end{aligned} \quad (1.22)$$

If the dimension or elements of Σ_i vary with i , then Ξ cannot be written as a direct product. The direct sum allows writing

$$\begin{aligned} \Xi &= \bigoplus_{i=1}^N \Sigma_i \\ &= \begin{bmatrix} \Sigma_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma_2 & \mathbf{0} & \vdots \\ \vdots & \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \Sigma_N \end{bmatrix}. \end{aligned} \quad (1.23)$$

1.5 RULES OF OPERATIONS

Unless otherwise specified, we assume \mathbf{A} and \mathbf{B} conform for the operations in question. Without additional knowledge of the matrices involved, the following are true.

Theorem 1.1 Some operations obey *commutative laws*:

$$\mathbf{A} + \mathbf{B} = \mathbf{B} + \mathbf{A} \quad (1.24)$$

$$\mathbf{A} - \mathbf{B} = \mathbf{B} - \mathbf{A} \quad (1.25)$$

$$\mathbf{A} \# \mathbf{B} = \mathbf{B} \# \mathbf{A} \quad (1.26)$$

$${}_a\mathbf{B} = \mathbf{B}{}_a. \quad (1.27)$$

It is important to recognize that $\mathbf{AB} \neq \mathbf{BA}$ and $\mathbf{A} \otimes \mathbf{B} \neq \mathbf{B} \otimes \mathbf{A}$, except in special cases.

Theorem 1.2 Some operations obey *associative laws*:

$$(\mathbf{A} + \mathbf{B}) + \mathbf{C} = \mathbf{A} + (\mathbf{B} + \mathbf{C}) \quad (1.28)$$

$$(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC}) \quad (1.29)$$

$$(\mathbf{A} \otimes \mathbf{B}) \otimes \mathbf{C} = \mathbf{A} \otimes (\mathbf{B} \otimes \mathbf{C}) \quad (1.30)$$

$$(\mathbf{aB}) \otimes (\mathbf{cD}) = \mathbf{ac}(\mathbf{B} \otimes \mathbf{D}) \quad (1.31)$$

$${}_a \otimes \mathbf{B} = \mathbf{B} \otimes {}_a = {}_a\mathbf{B}. \quad (1.32)$$

Theorem 1.3 Some operations obey *distributive laws*:

$$\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC} \quad (1.33)$$

$$\mathbf{A}(\mathbf{B} - \mathbf{C}) = \mathbf{AB} - \mathbf{AC} \quad (1.34)$$

$$(\mathbf{B} + \mathbf{C})\mathbf{D} = \mathbf{BD} + \mathbf{CD} \quad (1.35)$$

$$(\mathbf{B} - \mathbf{C})\mathbf{D} = \mathbf{BD} - \mathbf{CD} \quad (1.36)$$

$${}_a(\mathbf{B} + \mathbf{C}) = {}_a\mathbf{B} + {}_a\mathbf{C} = (\mathbf{B} + \mathbf{C}){}_a \quad (1.37)$$

$$(\mathbf{B} - \mathbf{C}){}_a = {}_a\mathbf{B} - {}_a\mathbf{C} = (\mathbf{B} - \mathbf{C}){}_a \quad (1.38)$$

$$(\mathbf{A} + \mathbf{B})' = \mathbf{A}' + \mathbf{B}' \quad (1.39)$$

$$(\mathbf{A} - \mathbf{B})' = \mathbf{A}' - \mathbf{B}'. \quad (1.40)$$

Theorem 1.4 The transpose has some special operational properties:

$$({}_a\mathbf{B})' = \mathbf{aB}' = \mathbf{B}'{}_a \quad (1.41)$$

$$(\mathbf{ABC} \dots)' = \dots \mathbf{C}'\mathbf{B}'\mathbf{A}' \quad (1.42)$$

$$(\mathbf{A} \otimes \mathbf{B} \otimes \mathbf{C} \otimes \dots)' = \mathbf{A}' \otimes \mathbf{B}' \otimes \mathbf{C}' \otimes \dots \quad (1.43)$$

$$\mathbf{ab}' = \mathbf{a} \otimes \mathbf{b}' = \mathbf{b}' \otimes \mathbf{a}. \quad (1.44)$$

Theorem 1.5 For conforming matrices,

$$\text{vec}(\mathbf{AB}) = (\mathbf{I} \otimes \mathbf{A})\text{vec}(\mathbf{B}) = (\mathbf{B}' \otimes \mathbf{I})\text{vec}(\mathbf{A}) \quad (1.45)$$

$$\text{vec}(\mathbf{ABC}') = (\mathbf{C}' \otimes \mathbf{A})\text{vec}(\mathbf{B}) \quad (1.46)$$

$$\text{vec}(\mathbf{ABCD}) = (\mathbf{I} \otimes \mathbf{A})(\mathbf{I} \otimes \mathbf{B})(\mathbf{I} \otimes \mathbf{C})\text{vec}(\mathbf{D}). \quad (1.47)$$

Theorem 1.6 (a) For any matrix pair, $\text{tr}(\mathbf{A} \otimes \mathbf{B}) = \text{tr}(\mathbf{A})\text{tr}(\mathbf{B})$.

(b) For conforming matrices, $\text{tr}(\mathbf{AB}) = \text{tr}(\mathbf{BA})$.

1.6 OTHER SPECIAL TYPES OF MATRICES

Definition 1.24 A matrix of the form $\mathbf{A}'\mathbf{A}$ is an *inner product* and \mathbf{AA}' is an *outer product*.

If $\mathbf{A} = [\mathbf{a}_1 \cdots \mathbf{a}_c]$ is $r \times c$, then $\mathbf{AA}' = \sum_{j=1}^c \mathbf{a}_j \mathbf{a}'_j$ is $r \times r$ and equals the sum of the c outer products of the c columns of \mathbf{A} . If $\mathbf{A} = \begin{bmatrix} \mathbf{A}_1 \\ \vdots \\ \mathbf{A}_r \end{bmatrix} = \begin{bmatrix} \mathbf{b}'_1 \\ \vdots \\ \mathbf{b}'_r \end{bmatrix} = \mathbf{B}'$ is $r \times c$, then $\mathbf{A}'\mathbf{A} = \sum_{k=1}^r \mathbf{b}_k \mathbf{b}'_k$ is the sum of r outer products of the columns of \mathbf{B} , which are the rows of \mathbf{A} . Both inner and outer products are always symmetric. Using concepts introduced later in the chapter, inner and outer products are always either positive definite (all eigenvalues real and positive) or positive semidefinite (all real eigenvalues, with some positive and some zero). Inner and outer products always have the same rank, which equals the rank of \mathbf{A} . They also have the same eigenvalues, except for some zeros if \mathbf{A} is not square.

Definition 1.25 A matrix is (columnwise) *orthogonal* if $\mathbf{A}'\mathbf{A}$ (the inner product) is diagonal, and a matrix is (rowwise) *orthogonal* if \mathbf{AA}' (the outer product) is diagonal. A matrix is (columnwise) *orthonormal* if $\mathbf{A}'\mathbf{A} = \mathbf{I}$, and a matrix is (rowwise) *orthonormal* if $\mathbf{AA}' = \mathbf{I}$. Two matrices are *biorthogonal* if $\mathbf{AB} = \mathbf{0}$.

In the preceding definition, neither \mathbf{A} nor \mathbf{B} need be square.

Definition 1.26 Any square matrix is described as *idempotent* if $\mathbf{A} = \mathbf{A}^2$.

Lemma 1.3 If \mathbf{A} is idempotent, then $\mathbf{I} - \mathbf{A}$ is also idempotent and $\mathbf{A}(\mathbf{I} - \mathbf{A}) = \mathbf{0}$.

Idempotent matrices play important roles in discovering properties of quadratic forms, especially independence.

1.7 QUADRATIC AND BILINEAR FORMS

Definition 1.27 (a) For square $\mathbf{A} = \mathbf{A}'$ and conforming \mathbf{x} , the expression $q = \mathbf{x}'\mathbf{A}\mathbf{x}$ is a *quadratic form* in \mathbf{x} .

(b) The expression $b = \mathbf{x}'_1 \mathbf{B} \mathbf{x}_2$, for \mathbf{B} not necessarily symmetric or square is a *bilinear form* in conforming vectors \mathbf{x}_1 and \mathbf{x}_2 .

If \mathbf{x} (2×1) is free to vary and $q_0 > 0$ is constant, then $q_0 = \mathbf{x}'\mathbf{A}\mathbf{x}$ is the equation of an ellipse. The result generalizes to higher dimensions. A quadratic form lies at the heart of the density of a vector Gaussian and consequently leads to ellipsoidal probability contours.

Lemma 1.4 If $q = \mathbf{x}'\mathbf{A}\mathbf{x}$, then without loss of generality \mathbf{A} may be assumed to be symmetric.

Proof. If $\mathbf{C} = (\mathbf{A} + \mathbf{A}')/2$ then $\mathbf{C} = \mathbf{C}'$ and

$$\begin{aligned} \mathbf{x}'\mathbf{C}\mathbf{x} &= \mathbf{x}'[(\mathbf{A} + \mathbf{A}')/2]\mathbf{x} \\ &= [\mathbf{x}'\mathbf{A}\mathbf{x} + (\mathbf{x}'\mathbf{A}\mathbf{x})']/2 \\ &= \mathbf{x}'\mathbf{A}\mathbf{x}. \end{aligned} \tag{1.48} \quad \square$$

Lemma 1.5 Without loss of generality, any bilinear form, $\mathbf{x}_1\mathbf{B}\mathbf{x}_2$, may be expressed as a quadratic form, $\mathbf{y}'\mathbf{D}\mathbf{y}$, with $\mathbf{y}' = [\mathbf{x}'_1 \ \mathbf{x}'_2]$ and $\mathbf{D} = \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ \mathbf{B}' & \mathbf{0} \end{bmatrix}/2$.

Proof. Here $\mathbf{D} = \mathbf{D}'$ and

$$\begin{aligned} \mathbf{y}'\mathbf{D}\mathbf{y} &= \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}' \left(\frac{1}{2} \begin{bmatrix} \mathbf{0} & \mathbf{B} \\ \mathbf{B}' & \mathbf{0} \end{bmatrix} \right) \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \\ &= [\mathbf{x}'_1 \ \mathbf{x}'_2] \begin{bmatrix} \mathbf{B}\mathbf{x}_2 \\ \mathbf{B}'\mathbf{x}_1 \end{bmatrix} / 2 \\ &= (\mathbf{x}'_1\mathbf{B}\mathbf{x}_2 + \mathbf{x}'_2\mathbf{B}'\mathbf{x}_1)/2 \\ &= \mathbf{x}'_1\mathbf{B}\mathbf{x}_2. \end{aligned} \tag{1.49} \quad \square$$

1.8 VECTOR SPACES AND RANK

Definition 1.28 A set of $n \times 1$ vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_p\}$ is *linearly dependent* if a set of scalar coefficients $\{a_1, \dots, a_p\}$, not all zero, exist such that

$$\sum_{i=1}^p a_i \mathbf{x}_i = a_1 \mathbf{x}_1 + a_2 \mathbf{x}_2 + \dots + a_p \mathbf{x}_p = \mathbf{0}. \tag{1.50}$$

If no such set of a_i exists then the set of \mathbf{x}_i is *linearly independent*. The single equation about $n \times 1$ vectors defines a set of n scalar equations.

Definition 1.29 (a) Any finite set of $n \times 1$ vectors generates a *vector space*, namely the (usually infinite) collection of all possible vectors created by any combination of multiplications of one vector by a constant, or the addition of two vectors.

(b) Any set of vectors which generate the particular set of vectors *spans* the vector space, and provides a *basis* for the space.

Definition 1.30 (a) The *rank* of the vector space equals the smallest possible number of linearly independent vectors which span the space. The rank of a set equals zero if and only if the only member of the set is $\mathbf{x}_i \equiv \mathbf{0}$. The rank of a set of p vectors, necessarily an integer, ranges from zero to p .
 (b) A set with rank p is *full rank*, while a set with rank strictly less than p is *less than full rank*.

Any two distinct vectors \mathbf{x}_1 and \mathbf{x}_2 are orthogonal if and only if $\mathbf{x}'_1 \mathbf{x}_2 = 0$. An orthogonal basis provides the most convenient form and has $\mathbf{x}'_j \mathbf{x}_{j'} = 0$ if $j \neq j'$. Spectral (eigenvalue) decomposition provides an orthonormal basis for any square and symmetric matrix, and some nonsymmetric square matrices. The singular value decomposition provides a convenient way for any matrix, symmetric or not, square or not. Both are discussed later in the chapter.

An $r \times c$ matrix, \mathbf{A} , can be thought of as a collection of c vectors, the columns, each $r \times 1$. Alternately, considering the columns of \mathbf{A}' allows describing the matrix as a collection of r vectors, the transposed rows, each $c \times 1$. The rank of a matrix may be found by decomposing it into its columns and treating them as a set of vectors:

$$\begin{aligned} \mathbf{A} &= \begin{bmatrix} 1 & 4 & 7 \\ 2 & 5 & 8 \\ 3 & 6 & 9 \\ 4 & 9 & 2 \end{bmatrix} \\ &= [\mathbf{a}_1 \quad \mathbf{a}_2 \quad \mathbf{a}_3] \end{aligned} \tag{1.51}$$

$$\Leftrightarrow \{\mathbf{a}_1, \mathbf{a}_2, \mathbf{a}_3\} = \left\{ \begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix}, \begin{bmatrix} 4 \\ 5 \\ 6 \\ 9 \end{bmatrix}, \begin{bmatrix} 7 \\ 8 \\ 9 \\ 2 \end{bmatrix} \right\}. \tag{1.52}$$

Transposing the matrix allows applying the same process to the rows. The resulting row rank always equals the column rank, which leads to the following.

Definition 1.31 The *rank of a matrix* equals the maximum number of linearly independent rows or columns, indicated $\text{rank}(\mathbf{A})$. An $r \times c$ matrix is *full rank* if $\text{rank}(\mathbf{A}) = \min(r, c)$ and *less than full rank* otherwise. The only matrix of rank zero is a matrix of all zeros, $\mathbf{0}_{n \times m}$.

It would be hard to overemphasize the importance of the concepts of vector space, span, basis, and orthogonal basis in the study of linear models. The concepts lie at the heart of many theorems, proofs and computational methods.

The same ideas provide crucial tools in understanding the logic of hypothesis tests, constrained models, and equivalence between models and parameterizations.

1.9 FINDING RANK

One method to find rank uses only elementary row operations or only elementary column operations to produce a canonical form, which is necessarily of equivalent rank. In particular, using (only) elementary row operations allows transforming a matrix to a triangular one (triangularize the matrix). The three elementary row (column) operations are: (1) multiplying a row (column) by a nonzero constant, (2) adding one row (column) to another, and (3) exchanging two rows (columns).

A second approach to finding rank uses the spectral or singular value decompositions described later in the chapter. The rank of any square symmetric matrix equals the number of nonzero eigenvalues (although the relationship may not hold for square but not symmetric matrices). However, *any* matrix has a singular value decomposition, with the number of nonzero singular values equal to the rank of the matrix.

The concepts of similarity and congruence (defined in Section 1.12) can be used to simplify the task. The following lemmas also prove useful.

Lemma 1.6 In general

$$0 \leq \text{rank}(\mathbf{A}) \leq \min(r, c) \tag{1.53}$$

$$\text{rank}(\mathbf{AB}) \leq \min[\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B})] \tag{1.54}$$

$$\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B}). \tag{1.55}$$

Lemma 1.7 (a) The rank of a diagonal matrix equals the number of nonzero diagonal elements.

(b) If $b \neq 0$, then

$$\text{rank}(\mathbf{A}) = \text{rank}(b\mathbf{A}). \tag{1.56}$$

The case of $b = -1$ provides an example: $\text{rank}(\mathbf{A}) = \text{rank}(-\mathbf{A})$.

(c) For any matrix

$$\text{rank}(\mathbf{A}) = \text{rank}(\mathbf{A}') = \text{rank}(\mathbf{A}'\mathbf{A}) = \text{rank}(\mathbf{A}\mathbf{A}'). \tag{1.57}$$

Lemma 1.8 Multiplying by a square full-rank matrix does not change rank. If \mathbf{A} is $r \times r$ of rank r , \mathbf{B} is $r \times c$ with $0 \leq \text{rank}(\mathbf{B}) \leq \min(r, c)$, and \mathbf{C} is $c \times c$ of rank c , then

$$\text{rank}(\mathbf{ABC}') = \text{rank}(\mathbf{AB}) = \text{rank}(\mathbf{BC}') = \text{rank}(\mathbf{B}). \tag{1.58}$$

Lemma 1.9 If \mathbf{A} is $r \times c$ of rank c and \mathbf{B} is $c \times d$ with $\text{rank}(\mathbf{B}) = c$, then

$$\text{rank}(\mathbf{A}\mathbf{B}) = \text{rank}(\mathbf{A}). \quad (1.59)$$

Lemma 1.10 For any \mathbf{A} and \mathbf{B}

$$\text{rank}(\mathbf{A} \otimes \mathbf{B}) = \text{rank}(\mathbf{A})\text{rank}(\mathbf{B}) \quad (1.60)$$

and for any square \mathbf{A} and any square \mathbf{B}

$$\text{rank}(\mathbf{A} \oplus \mathbf{B}) = \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B}). \quad (1.61)$$

1.10 DETERMINANTS

Definition 1.32 The *determinant*, a scalar, is indicated $|\mathbf{A}|$ and is defined only for a square matrix. Also $|\mathbf{A}'| = |\mathbf{A}|$.

For any 2×2 matrix

$$\left| \begin{bmatrix} a & b \\ c & d \end{bmatrix} \right| = ad - bc. \quad (1.62)$$

The determinant of a diagonal or triangular matrix equals the product of the diagonal values. For 3×3 matrices

$$\left| \begin{bmatrix} a & 0 & 0 \\ 0 & b & 0 \\ 0 & 0 & c \end{bmatrix} \right| = \left| \begin{bmatrix} a & 0 & 0 \\ x & b & 0 \\ y & z & c \end{bmatrix} \right| = \left| \begin{bmatrix} a & x & y \\ 0 & b & z \\ 0 & 0 & c \end{bmatrix} \right| = abc. \quad (1.63)$$

For any general 3×3 matrix

$$\left| \begin{bmatrix} a & b & c \\ d & e & f \\ g & h & i \end{bmatrix} \right| = aei + bfg + chd - ceg - afh - bdi. \quad (1.64)$$

In general the determinant equals the sum of products with alternating sign of elements of the matrix. The determinant equals the product of the eigenvalues (discussed later in the present chapter), which provides the most useful interpretation of the determinant for statistical analysis. Determinants have many useful properties related to rank.

Lemma 1.11 For \mathbf{A} ($n \times n$),

$$\begin{aligned} |\mathbf{A}| = 0 &\Leftrightarrow \text{rank}(\mathbf{A}) < n \\ &\Leftrightarrow \mathbf{A} \text{ less than full rank} \\ &\Leftrightarrow \mathbf{A}^{-1} \text{ does not exist,} \end{aligned}$$

while

$$\begin{aligned}
 |\mathbf{A}| \neq 0 &\Leftrightarrow \text{rank}(\mathbf{A}) = n \\
 &\Leftrightarrow \mathbf{A} \text{ full rank} \\
 &\Leftrightarrow \mathbf{A}^{-1} \text{ exists,} \\
 &\Leftrightarrow \text{all columns (rows) of } \mathbf{A} \text{ linearly independent.}
 \end{aligned}$$

Definition 1.33 A less-than-full-rank square matrix may be described as *singular* and a full-rank square matrix as *nonsingular*.

Lemma 1.12 (a) $|\mathbf{A}| = |\mathbf{A}'|$.

(b) If \mathbf{A} is $m \times m$ and \mathbf{B} is $n \times n$ then $|\mathbf{A} \oplus \mathbf{B}| = |\mathbf{A}||\mathbf{B}|$ and

(c) $|\mathbf{A} \otimes \mathbf{B}| = |\mathbf{A}|^m |\mathbf{B}|^n$.

(d) If $m = n$, then $|\mathbf{AB}| = |\mathbf{A}||\mathbf{B}| = |\mathbf{BA}|$.

1.11 THE INVERSE AND GENERALIZED INVERSE

Theorem 1.7 (a) An $n \times n$ matrix \mathbf{A} of full rank has $\text{rank}(\mathbf{A}) = n$, and there exists a unique matrix \mathbf{A}^{-1} , called the *inverse* of \mathbf{A} , such that

$$\mathbf{A}^{-1}\mathbf{A} = \mathbf{A}\mathbf{A}^{-1} = \mathbf{I}. \tag{1.65}$$

(b) If $n = 1$, then \mathbf{A} is a scalar, and $\mathbf{A}^{-1} = 1/a$ exists if and only if $a \neq 0$.

(c) The inverse of a full-rank diagonal matrix equals the diagonal matrix of reciprocals of the diagonal elements.

Lemma 1.13 For square, conforming (same-size) and full-rank matrices

$$(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}. \tag{1.66}$$

More generally, for any finite set of full-rank (and same size) matrices

$$(\mathbf{ABCD}\dots)^{-1} = \dots\mathbf{D}^{-1}\mathbf{C}^{-1}\mathbf{B}^{-1}\mathbf{A}^{-1}. \tag{1.67}$$

The inverse has some symmetries. Not surprisingly, a symmetric matrix has a symmetric inverse. When they exist, the inverse of the transpose equals the transpose of the inverse and hence may be written unequivocally as \mathbf{A}^{-t} :

$$\mathbf{A}^{-t} = (\mathbf{A}')^{-1} = (\mathbf{A}^{-1})'. \tag{1.68}$$

Furthermore, if \mathbf{A}^{-1} exists, then $|\mathbf{A}^{-1}| = (|\mathbf{A}|)^{-1}$.

Lemma 1.14 (a) If $\mathbf{V} = \mathbf{R} + \mathbf{ZGZ}'$ and \mathbf{V} , \mathbf{R} , and \mathbf{G} are all square and full rank, then $\mathbf{V}^{-1} = \mathbf{R}^{-1}(\mathbf{R} - \mathbf{ZBZ}')\mathbf{R}^{-1}$, with $\mathbf{B} = (\mathbf{G}^{-1} + \mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z})^{-1} = \mathbf{G}(\mathbf{G}^{-1} - \mathbf{Z}'\mathbf{V}^{-1}\mathbf{Z})\mathbf{G}$.

(b) If $\mathbf{V} = \mathbf{R} + g\mathbf{ZZ}'$ and \mathbf{V} , \mathbf{R} , and g are all square and full rank, then

$(\mathbf{R} + g\mathbf{Z}\mathbf{Z}')^{-1} = \mathbf{R}^{-1}(\mathbf{R} - b\mathbf{Z}\mathbf{Z}')\mathbf{R}^{-1}$, with $b = (g^{-1} + \mathbf{Z}'\mathbf{R}^{-1}\mathbf{Z})^{-1}$.

(c) For \mathbf{A} ($r \times c$) and \mathbf{B} ($c \times r$), if $\text{rank}(\mathbf{I} + \mathbf{A}\mathbf{B}) = r$, then $(\mathbf{I} + \mathbf{A}\mathbf{B})^{-1} = \mathbf{I} - \mathbf{A}(\mathbf{I} + \mathbf{B}\mathbf{A})^{-1}\mathbf{B}$.

(d) Applying (c) recursively gives $(\mathbf{I}_r - \mathbf{C})^{-1} = \mathbf{I} + \sum_{j=1}^{\infty} \mathbf{C}^j$ (when it exists).

Theorem 1.8 For any matrix \mathbf{A} , square or not, a nonunique *generalized inverse* \mathbf{A}^- always exists such that

$$\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}. \quad (1.69)$$

One particularly inconvenient property of the one-condition generalized inverse is that \mathbf{A}^- need not be symmetric even though \mathbf{A} is symmetric.

Theorem 1.9 For any matrix \mathbf{A} , square or not, the unique *Moore-Penrose (generalized) inverse* \mathbf{A}^+ always exists and satisfies four conditions:

$$1. \quad \mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A} \quad (1.70)$$

$$2. \quad \mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+ \quad (1.71)$$

$$3. \quad (\mathbf{A}^+\mathbf{A})' = \mathbf{A}^+\mathbf{A} \quad (1.72)$$

$$4. \quad (\mathbf{A}\mathbf{A}^+)' = \mathbf{A}\mathbf{A}^+. \quad (1.73)$$

If \mathbf{A} is symmetric, then \mathbf{A}^+ is always symmetric. Matrices meeting only subsets of conditions, such as 1 and 2, (a two-condition inverse), or conditions 1, 2, and 3 (a three-condition inverse) have also been studied. Although a few proofs in linear models require a two- or three-condition inverse, we nearly always restrict attention to either one- or four-condition inverses.

Both the Moore-Penrose (four-condition) and any one-condition generalized inverse always coincide with the regular inverse for a square and full-rank matrix. Furthermore any one-condition inverse always coincides with the four-condition inverses for a nonsquare and full-rank matrix.

For a wide variety of linear models with a less-than-full-rank design matrix \mathbf{X} , the form $(\mathbf{X}'\mathbf{X})^-$ occurs often in expressions for estimators and distribution parameters. Very conveniently, Theorem 1.15 in the next section eliminates the need to distinguish between $(\mathbf{X}'\mathbf{X})^-$ and $(\mathbf{X}'\mathbf{X})^+$ in a wide range of linear model applications.

Theorem 1.10 (a) The Moore-Penrose inverse of a less-than-full-rank diagonal matrix is the diagonal matrix with reciprocals of the nonzero elements on the diagonal in the same locations as the nonzero elements and zero elsewhere.

(b) The Moore-Penrose inverse of the transpose equals the transpose of the generalized inverse and hence may be written unequivocally as \mathbf{A}^{+t} :

$$\mathbf{A}^{+t} = (\mathbf{A}')^+ = (\mathbf{A}^+)'. \quad (1.74)$$

(c) If $r \times c$ \mathbf{A} has $\text{rank}(\mathbf{A}) = r \leq c$, then \mathbf{A}^+ has dimensions $c \times r$ and

$$\mathbf{A}^+ = \mathbf{A}'(\mathbf{A}\mathbf{A}')^{-1}. \quad (1.75)$$

(d) If $r \times c$ \mathbf{A} has $r \geq c = \text{rank}(\mathbf{A})$, then \mathbf{A}^+ has dimensions $c \times r$ and

$$\mathbf{A}^+ = (\mathbf{A}'\mathbf{A})^{-1}\mathbf{A}'. \quad (1.76)$$

In contrast to the regular inverse, $(\mathbf{A}\mathbf{B})^+$ may or may not equal $\mathbf{B}^+\mathbf{A}^+$. The next lemma lists two side conditions which suffice to ensure that the result holds.

Lemma 1.15 (a) If \mathbf{A} is $r \times c$, \mathbf{P} is $r \times r$, and \mathbf{Q} is $c \times c$, with $\mathbf{P}'\mathbf{P} = \mathbf{I}_r$ and $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_c$, then

$$(\mathbf{P}\mathbf{A}\mathbf{Q}')^+ = \mathbf{P}\mathbf{A}^+\mathbf{Q}'. \quad (1.77)$$

(b) If \mathbf{A} is $r \times c$ of rank c and \mathbf{B} is $c \times d$ of rank c , then

$$(\mathbf{A}\mathbf{B})^+ = \mathbf{B}^+\mathbf{A}^+. \quad (1.78)$$

Definition 1.34 For known \mathbf{A} ($r \times c$) and known \mathbf{b} ($r \times 1$), system $\mathbf{A}\mathbf{x} = \mathbf{b}$ is *consistent* whenever any linear relationships existing among the rows of \mathbf{A} also exist among the corresponding rows of \mathbf{b} ($r \times 1$). Equivalently $\mathbf{c}'\mathbf{A} = \mathbf{0} \Rightarrow \mathbf{c}'\mathbf{b} = 0$, \Leftrightarrow the system has one or more solutions, \mathbf{x}_0 .

By the definition, a system of equations is consistent if at least one solution set exists. More than one or even infinitely many solutions may exist. With \mathbf{A} and \mathbf{b} known constants and \mathbf{x} unknown, the equation $\mathbf{A}\mathbf{x} = \mathbf{b}$ defines a system of equations (one per row of \mathbf{A} and \mathbf{b}).

Lemma 1.16 For known \mathbf{A} ($r \times c$) and known \mathbf{b} ($r \times 1$), if the system $\mathbf{A}\mathbf{x} = \mathbf{b}$ is consistent, then the following all hold.

(a) If $r = c$ and $\text{rank}(\mathbf{A}) = r$, then the system has a unique solution given by $\mathbf{x} = \mathbf{A}^{-1}\mathbf{b}$.

(b) If $r \neq c$ or $\text{rank}(\mathbf{A}) \neq r$ or both and \mathbf{A}^- is any generalized inverse of \mathbf{A} (i.e., $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A}$), then $\mathbf{x} = \mathbf{A}^-\mathbf{b}$ is one solution of infinitely many.

Proof of (b). $\mathbf{A}\mathbf{A}^-\mathbf{A} = \mathbf{A} \Rightarrow (\mathbf{A}\mathbf{A}^-\mathbf{A})\mathbf{x} = \mathbf{b} \Rightarrow \mathbf{A}\mathbf{A}^-(\mathbf{A}\mathbf{x}) = \mathbf{b} \Rightarrow \mathbf{A}\mathbf{A}^-\mathbf{b} = \mathbf{b}$. □

Lemma 1.17 If \mathbf{A} ($r \times c$) has $(\mathbf{A}^-)_1$ as a particular one-condition generalized inverse and \mathbf{B} is $c \times r$, then

$$(\mathbf{A}^-)_2 = (\mathbf{A}^-)_1 + \mathbf{B} - (\mathbf{A}^-)_1\mathbf{A}\mathbf{B}\mathbf{A}(\mathbf{A}^-)_1 \quad (1.79)$$

is also a one-condition generalized inverse for \mathbf{A} .

Corollary 1.17 For \mathbf{A} ($r \times c$) and any \mathbf{A}^- ($c \times r$), \mathbf{B} and \mathbf{C} exist such that

$$\mathbf{A}^- = \mathbf{A}^+ + \mathbf{B} - \mathbf{A}^+ \mathbf{A} \mathbf{B} \mathbf{A} \mathbf{A}^+ \quad (1.80)$$

$$\mathbf{A}^+ = \mathbf{A}^- + \mathbf{C} - \mathbf{A}^- \mathbf{A} \mathbf{C} \mathbf{A} \mathbf{A}^- . \quad (1.81)$$

Lemma 1.18 If \mathbf{A}^{-1} exists, then

$$\text{rank}(\mathbf{A}^{-1}) = \text{rank}(\mathbf{A}) . \quad (1.82)$$

For \mathbf{A} of any dimension and any rank

$$\text{rank}(\mathbf{A}^+) = \text{rank}(\mathbf{A} \mathbf{A}^+) = \text{rank}(\mathbf{A}^+ \mathbf{A}) = \text{rank}(\mathbf{A}) \quad (1.83)$$

and

$$\text{rank}(\mathbf{A}^-) \geq \text{rank}(\mathbf{A} \mathbf{A}^-) = \text{rank}(\mathbf{A}^- \mathbf{A}) = \text{rank}(\mathbf{A}) . \quad (1.84)$$

Lemma 1.19 If \mathbf{A} is $r \times c$, then $\mathbf{A}^- \mathbf{A}$ is $c \times c$, $\mathbf{A} \mathbf{A}^-$ is $r \times r$, and both are idempotent. The same properties hold if \mathbf{A}^+ replaces \mathbf{A}^- .

Proof. $\mathbf{A}^- \mathbf{A} = \mathbf{A}^- \mathbf{A} \mathbf{A}^- \mathbf{A} = (\mathbf{A}^- \mathbf{A})^2$ and $\mathbf{A} \mathbf{A}^- = \mathbf{A} \mathbf{A}^- \mathbf{A} \mathbf{A}^- = (\mathbf{A} \mathbf{A}^-)^2$. \square

Lemma 1.20 For any \mathbf{A} and \mathbf{B}

$$(\mathbf{A} \otimes \mathbf{B})^- = \mathbf{A}^- \otimes \mathbf{B}^- \quad (1.85)$$

$$(\mathbf{A} \otimes \mathbf{B})^+ = \mathbf{A}^+ \otimes \mathbf{B}^+ , \quad (1.86)$$

and for any square \mathbf{A} and square \mathbf{B}

$$(\mathbf{A} \oplus \mathbf{B})^- = \mathbf{A}^- \oplus \mathbf{B}^- \quad (1.87)$$

$$(\mathbf{A} \oplus \mathbf{B})^+ = \mathbf{A}^+ \oplus \mathbf{B}^+ . \quad (1.88)$$

If square \mathbf{A} and square \mathbf{B} are both full rank, then

$$(\mathbf{A} \oplus \mathbf{B})^{-1} = \mathbf{A}^{-1} \oplus \mathbf{B}^{-1} \quad (1.89)$$

$$(\mathbf{A} \otimes \mathbf{B})^{-1} = \mathbf{A}^{-1} \otimes \mathbf{B}^{-1} . \quad (1.90)$$

1.12 EIGENANALYSIS (SPECTRAL DECOMPOSITION)

Eigenanalysis is only defined for square \mathbf{A} . Nearly all interest in decomposing matrices in statistics lies with symmetric matrices. The symmetry (in statistical applications) arises because the matrices of interest are inner products or outer products. However, for the moment, we consider any square matrix, $n \times n$, symmetric or not.

Definition 1.35 A *right eigenvector* of \mathbf{A} is an $n \times 1$ vector, $\mathbf{x} \neq \mathbf{0}$, such that

$$\mathbf{Ax} = \lambda \mathbf{x}, \tag{1.91}$$

with λ the *eigenvalue* corresponding to \mathbf{x} . A *left eigenvector* is an $n \times 1$ vector \mathbf{y} such that $\mathbf{y}'\mathbf{A} = \lambda\mathbf{y}'$. A left eigenvector of \mathbf{A} is a right eigenvector of \mathbf{A}' : $\mathbf{A}'\mathbf{y} = \lambda\mathbf{y}$. The set of eigenvalues is sometimes referred to as the *spectrum* of the matrix.

Eigen means “characteristic” in German, which leads to the alternative descriptions as *characteristic values* and *characteristic vectors*. The eigenvector \mathbf{x} has the special property of projecting \mathbf{A} back into itself (\mathbf{x}) times a constant. The self-replicating feature leads to thinking of eigenvectors as matrix DNA.

Statisticians most often apply eigenanalysis to symmetric matrices. Right and left eigenvectors coincide for a symmetric matrix but usually do not for a nonsymmetric matrix and may not even exist. One important exception occurs in multivariate linear models for which the test statistics and associated contrasts correspond to eigenvalues and eigenvectors of a nonsymmetric matrix. For computational purposes, although not for scientific interpretation, the task can be expressed in terms of a symmetric matrix.

Theorem 1.11 (a) The roots of the *characteristic equation*, $|\mathbf{A} - \lambda\mathbf{I}| = 0$, equal the eigenvalues (characteristic values).

(b) Also

$$\mathbf{Ax} = \lambda \mathbf{x} \tag{1.92}$$

\Leftrightarrow

$$(\mathbf{A} - \lambda\mathbf{I})\mathbf{x} = \mathbf{0} \tag{1.93}$$

\Leftrightarrow

$$|\mathbf{A} - \lambda\mathbf{I}| = 0. \tag{1.94}$$

(c) The characteristic equation of an $n \times n$ matrix equals a polynomial in λ of order n .

For a 2×2

$$|\mathbf{A} - \lambda\mathbf{I}| = \left| \begin{bmatrix} a & b \\ c & d \end{bmatrix} - \lambda \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \right| \tag{1.95}$$

implies

$$\begin{aligned} 0 &= \left| \begin{bmatrix} a & b \\ c & d \end{bmatrix} - \begin{bmatrix} \lambda & 0 \\ 0 & \lambda \end{bmatrix} \right| = \begin{vmatrix} a - \lambda & b \\ c & d - \lambda \end{vmatrix} \\ &= (a - \lambda)(d - \lambda) - bc \\ &= \lambda^2 - \lambda(a + d) + ad - bc. \end{aligned} \tag{1.96}$$

Of course the quadratic formula allows solving for λ here and gives two values.

The properties of eigenvalues and eigenvectors vary substantially, with many awkward possibilities. As for the roots of any polynomial in real numbers, the eigenvalues of a real matrix may be imaginary. Furthermore, even though an $n \times n$ matrix always has n eigenvalues, not all $n \times n$ matrices have a complete set of n distinct eigenvectors. The number of distinct eigenvalues ranges from 1 to n .

Although distinct eigenvalues are unique, eigenvectors (when they exist) are not completely unique. If \mathbf{x} is an eigenvector of \mathbf{A} , then $c\mathbf{x}$ is also, for $c \neq 0$, because $\mathbf{A}\mathbf{x} = \lambda\mathbf{x}$ implies $\mathbf{A}c\mathbf{x} = \lambda c\mathbf{x}$. By strong convention, \mathbf{x} is usually scaled such that $\mathbf{x}'\mathbf{x} = 1$, described as *normalized to unit length*. However, a normalized eigenvector is still not unique because both \mathbf{x} and $-\mathbf{x}$ are normalized eigenvectors for \mathbf{A} . Geometrically, reversing the sign of the vector corresponds to a reflection about an axis. Although the sign ambiguity affects some computations, the subtlety can usually be ignored in discussing “the” (normalized) eigenvectors.

Definition 1.36 (a) *Algebraic multiplicity* is the number of times a particular distinct eigenvalue, λ_j , occurs as a root of the characteristic equation.
(b) The *geometric multiplicity* of λ_j equals the number of distinct eigenvectors associated with λ_j .
(c) A *simple* matrix has geometric multiplicity equal to algebraic multiplicity.

It follows from the definitions that geometric multiplicity must be less than or equal to algebraic multiplicity. The definitions have the important implication that any simple matrix (necessarily square but not necessarily symmetric) has a decomposition in terms of the diagonal matrix of eigenvalues.

Theorem 1.12 (a) All symmetric matrices are simple.

- (b)** A real and symmetric matrix has real eigenvalues and eigenvectors.
- (c)** Any real and symmetric matrix ($n \times n$) has a *spectral* decomposition,

$$\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{g}(\lambda)\mathbf{V}', \tag{1.97}$$

with \mathbf{V} a set of orthogonal eigenvectors and λ the corresponding eigenvalues in the same order (usually sorted from largest to smallest).

(d) By convention, and without loss of generality, $\mathbf{V}'\mathbf{V} = \mathbf{V}\mathbf{V}' = \mathbf{I}_n$, giving orthonormal \mathbf{V} and eigenvectors ($\mathbf{v}'_j\mathbf{v}_j = 1, \mathbf{v}'_j\mathbf{v}_k = 0$ for $j \neq k$).

(e) If \mathbf{D} ($n \times n$) is a diagonal matrix with diagonal elements freely chosen from $\{+1, -1\}$, without loss of generality \mathbf{V} may be taken to be $\mathbf{V}\mathbf{D}$.

The columns of $\mathbf{V} = [\mathbf{v}_1 \mathbf{v}_2 \cdots \mathbf{v}_n]$ are “the” normalized eigenvectors of \mathbf{A} , corresponding to the eigenvalues. The fact that $\mathbf{D}\mathbf{D} = \mathbf{I}$ gives $\mathbf{D}\mathbf{D}\mathbf{g}(\lambda)\mathbf{D} = \mathbf{D}\mathbf{g}(\lambda)$ and $\mathbf{A} = (\mathbf{V}\mathbf{D})\mathbf{D}\mathbf{g}(\lambda)(\mathbf{V}\mathbf{D})'$. Eigenvalues and eigenvectors *must* be in a corresponding order. Also \mathbf{V} always has full rank and $\mathbf{V}^{-1} = \mathbf{V}'$ (true for square, orthonormal matrices). Without loss of generality, but for convenience and by strong convention, \mathbf{V} is scaled to unit length ($\mathbf{v}'_j\mathbf{v}_j = 1$). A simple example is

$$\begin{aligned}
 \mathbf{A} &= \begin{bmatrix} 1.5 & -0.5 \\ -0.5 & 1.5 \end{bmatrix} \\
 &= \mathbf{V}\text{Dg}(\boldsymbol{\lambda})\mathbf{V}' \\
 &= \left(\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} / \sqrt{2} \right) \begin{bmatrix} 2 & 0 \\ 0 & 1 \end{bmatrix} \left(\begin{bmatrix} 1 & -1 \\ 1 & 1 \end{bmatrix} / \sqrt{2} \right)', \tag{1.98}
 \end{aligned}$$

with $\boldsymbol{\lambda} = [2 \ 1]'$, $\mathbf{V} = [\mathbf{v}_1 \ \mathbf{v}_2]$.

Despite the many awkward possibilities, statistical applications of eigenanalysis nearly always involve well behaved matrices such as $\mathbf{X}'\mathbf{X}$ or $\mathbf{X}\mathbf{X}'$. Such inner and outer products are always symmetric and therefore simple with real eigenvalues. Furthermore, they never have negative eigenvalues (all eigenvalues are positive or zero; nonnegative definite). Sums of squares, covariance, and correlation matrices can all be expressed as inner or outer products. In the study of linear models, the few computations requiring eigenanalysis of nonsymmetric matrices can be expressed in terms of closely related symmetric matrices.

Definition 1.37 (a) If \mathbf{A} , \mathbf{B} , and \mathbf{T} are $n \times n$ with \mathbf{T} full rank, then \mathbf{A} and \mathbf{B} are said to be *similar* if and only if $\mathbf{B} = \mathbf{T}\mathbf{A}\mathbf{T}^{-1}$.
(b) Matrices \mathbf{A} and \mathbf{B} are *congruent* if and only if $\mathbf{B} = \mathbf{T}\mathbf{A}\mathbf{T}'$.

Lemma 1.21 (a) Similar matrices have all the same eigenvalues.

- (b)** Any matrix similar to a diagonal matrix has rank equal to the number of nonzero eigenvalues.
- (c)** Every symmetric matrix is similar to a diagonal matrix, namely a diagonal matrix of the eigenvalues of \mathbf{A} , and always has rank equal to the number of nonzero eigenvalues.
- (d)** Geometrically, the eigenvectors corresponding to the nonzero eigenvalues of a simple matrix (including all symmetric matrices) span and provide an orthogonal basis for the full rank subspace spanned by \mathbf{A} .
- (e)** Sylvester's law of inertia (Lancaster, 1969, p90) guarantees that congruent matrices have the same number of positive, negative and zero eigenvalues (but not necessarily the same values).

In contrast to the nice properties in the last lemma, the rank of a square and nonsymmetric matrix may not equal the number of nonzero eigenvalues. Although the matrix $\mathbf{A} = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ has rank 1, both eigenvalues are zero.

Lemma 1.22 For any square matrix, the trace and determinant equal functions of the eigenvalues:

$$\text{tr}(\mathbf{A}) = \sum_{j=1}^n \lambda_j \tag{1.99}$$

and

$$|\mathbf{A}| = \prod_{j=1}^n \lambda_j. \quad (1.100)$$

Lemma 1.23 (a) $|\mathbf{A}| = 0 \Leftrightarrow$ at least one eigenvalue equals zero $\Leftrightarrow \mathbf{A}$ is less than full rank; $|\mathbf{A}| \neq 0 \Leftrightarrow$ no eigenvalue equals zero $\Leftrightarrow \mathbf{A}$ is full rank.

(b) The eigenvalues of a (square) diagonal or triangular matrix are the diagonal elements.

(c) For square and full-rank \mathbf{B} and square \mathbf{A} of the same dimension, the eigenvalues of \mathbf{BAB}^{-1} are the eigenvalues of \mathbf{A} .

(d) The eigenvalues of a (square) orthonormal matrix are ± 1 .

Definition 1.38 (a) Square \mathbf{A} is *positive definite* if $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$ for any conforming finite \mathbf{x} .

(b) Square \mathbf{A} is *positive semidefinite* if $\mathbf{x}'\mathbf{A}\mathbf{x} \geq 0$ and $\mathbf{x}'\mathbf{A}\mathbf{x} = 0$ for at least one $\mathbf{x} \neq \mathbf{0}$.

(c) Square \mathbf{A} is *negative definite* if $\mathbf{x}'\mathbf{A}\mathbf{x} < 0$.

(d) Square \mathbf{A} is *negative semidefinite* if $\mathbf{x}'\mathbf{A}\mathbf{x} \leq 0$ and $\mathbf{x}'\mathbf{A}\mathbf{x} = 0$ for at least one $\mathbf{x} \neq \mathbf{0}$.

(e) A *nonnegative definite* matrix is either positive definite or positive semidefinite.

(f) A *nonpositive definite* matrix is either negative definite or negative semidefinite.

A symmetric matrix \mathbf{A} with eigenvalues $\{\lambda_j\}$ and $\min(\lambda_j) > 0$ is always *positive definite*. If $\min(\lambda_j) = 0$, then \mathbf{A} is *positive semidefinite*. Similarly, \mathbf{A} is *negative definite* if $\max(\lambda_j) < 0 \Leftrightarrow \mathbf{x}'\mathbf{A}\mathbf{x} < 0$ and *negative semidefinite* if $\max(\lambda_j) = 0$. Positive definite and negative definite matrices are full rank.

Inner and outer products, such as $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}\mathbf{X}'$, are symmetric and therefore simple with a spectral decomposition. They are also necessarily positive definite or positive semidefinite (nonnegative definite).

Lemma 1.24 (a) For symmetric $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{g}(\lambda)\mathbf{V}'$ of full rank

$$\mathbf{A}^{-1} = \mathbf{V}\mathbf{D}\mathbf{g}(\lambda)^{-1}\mathbf{V}'. \quad (1.101)$$

(b) For symmetric \mathbf{A} of any rank

$$\mathbf{A}^+ = \mathbf{V}[\mathbf{D}\mathbf{g}(\lambda)]^+\mathbf{V}'. \quad (1.102)$$

(c) For symmetric $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{g}(\lambda)\mathbf{V}'$ of any rank and finite integer $p > 0$

$$\mathbf{A}^p = \mathbf{V}[\mathbf{D}\mathbf{g}(\lambda)]^p\mathbf{V}'. \quad (1.103)$$

(d) For symmetric $\mathbf{A} = \mathbf{V}\mathbf{D}\mathbf{g}(\lambda)\mathbf{V}'$ of full rank and finite integer $p > 0$

$$\mathbf{A}^{-p} = (\mathbf{A}^{-1})^p = (\mathbf{A}^p)^{-1} = \mathbf{V}[\text{Dg}(\boldsymbol{\lambda})]^{-p}\mathbf{V}'. \quad (1.104)$$

(e) For symmetric $\mathbf{A} = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})\mathbf{V}'$ of any rank and finite integer $p > 0$

$$(\mathbf{A}^+)^p = (\mathbf{A}^p)^+ = \mathbf{V}\{[\text{Dg}(\boldsymbol{\lambda})]^p\}^+\mathbf{V}'. \quad (1.105)$$

The lemma illustrates the extent to which the spectral decomposition characterizes the matrix. Merely computing the reciprocals of the eigenvalues allows computing the inverse or generalized inverse of a symmetric matrix. Although we do not pursue the idea here, the concept of a matrix function of a symmetric matrix is well defined. If a function $f(\lambda_j)$ has a valid Taylor series expansion for the eigenvalues of $\mathbf{A} = \mathbf{A}'$, namely $\{\lambda_j\}$, then $f(\mathbf{A}) = \mathbf{V}[\text{Dg}(\{f(\lambda_j)\})]\mathbf{V}'$ is a well-defined matrix.

Lemma 1.25 (a) For $\mathbf{V} = [\mathbf{v}_j \cdots \mathbf{v}_n]$ with $\mathbf{V}'\mathbf{V} = \mathbf{V}\mathbf{V}' = \mathbf{I}_n$, the $n \times n$ and symmetric matrix \mathbf{A} with $\text{rank}(\mathbf{A}) = n_1 \leq n$ has n_1 nonzero eigenvalues and spectral decomposition $\mathbf{A} = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})\mathbf{V}'$.

(b) The corresponding *constituent matrix decomposition* is

$$\mathbf{A} = \sum_{j=1}^{n_1} \lambda_j \mathbf{v}_j \mathbf{v}_j' = \sum_{j=1}^{n_1} \lambda_j \mathbf{v}_j \mathbf{v}_j. \quad (1.106)$$

(c) Constituent matrix $\mathbf{G}_j = \mathbf{v}_j \mathbf{v}_j'$ is symmetric and idempotent, which gives $\text{rank}(\mathbf{v}_j \mathbf{v}_j') = \text{tr}(\mathbf{v}_j \mathbf{v}_j') = \text{tr}(\mathbf{v}_j' \mathbf{v}_j) = 1$.

(d) Aggregating the eigenvalues and corresponding eigenvectors into two mutually exclusive and together exhaustive groups of sizes n_1 and n_0 , with $n = n_1 + n_0$, allows writing

$$\begin{aligned} \mathbf{A} &= [\mathbf{V}_1 \ \mathbf{V}_0] \begin{bmatrix} \text{Dg}(\boldsymbol{\lambda}_1) & \mathbf{0} \\ \mathbf{0} & \text{Dg}(\boldsymbol{\lambda}_0) \end{bmatrix} \begin{bmatrix} \mathbf{V}_1' \\ \mathbf{V}_0' \end{bmatrix} \\ &= \mathbf{V}_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{V}_1' + \mathbf{V}_0 \text{Dg}(\boldsymbol{\lambda}_0) \mathbf{V}_0'. \end{aligned} \quad (1.107)$$

(e) If $\boldsymbol{\lambda}_1$ contains the n_1 nonzero eigenvalues then $\boldsymbol{\lambda}_0 = \mathbf{0}$ and $\mathbf{A} = \mathbf{V}_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{V}_1'$, with \mathbf{V}_1 $n \times n_1$ ($\mathbf{V}_1' \mathbf{V}_1 = \mathbf{I}_{n_1}$), while $\text{Dg}(\boldsymbol{\lambda}_1)$ is $n_1 \times n_1$ of full rank.

(f) The decomposition in part (d), in terms of two groups of eigenvalues and vectors, generalizes to three or more groups.

Lemma 1.26 Any $n \times n$ symmetric and idempotent matrix \mathbf{A} of rank $n_1 \leq n$ has

(a) n_1 eigenvalues of 1, (b) all remaining $n - n_1$ eigenvalues of 0, and (c) $\text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A})$. (d) The eigenvectors $\mathbf{V} = [\mathbf{V}_1 \ \mathbf{V}_0]$ may be arranged and scaled such that $\mathbf{V}_1' \mathbf{V}_1 = \mathbf{I}_{n_1}$, $\mathbf{V}_0' \mathbf{V}_0 = \mathbf{I}_{n-n_1}$, $\mathbf{V}_1' \mathbf{V}_0 = \mathbf{0}$, and

$$\mathbf{A} = [\mathbf{V}_1 \ \mathbf{V}_0] \begin{bmatrix} \mathbf{I}_{n_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{n-n_1} \end{bmatrix} \begin{bmatrix} \mathbf{V}_1' \\ \mathbf{V}_0' \end{bmatrix} = \mathbf{V}_1 \mathbf{V}_1'. \quad (1.108)$$

Lemma 1.27 (a) If $m \times 1$ \mathbf{a} contains the m eigenvalues of $m \times m$ \mathbf{A} , while $n \times 1$ \mathbf{b} contains the n eigenvalues of $n \times n$ \mathbf{B} , then the eigenvalues of $\mathbf{C} = \mathbf{A} \oplus \mathbf{B}$ are the $m + n$ elements of $\mathbf{c} = [\mathbf{a}' \ \mathbf{b}']'$.

(b) The eigenvalues of $\mathbf{D} = \mathbf{A} \otimes \mathbf{B}$ are the mn elements of $\mathbf{d} = \mathbf{a} \otimes \mathbf{b}$. Neither \mathbf{c} nor \mathbf{d} is necessarily sorted by size.

Many regression diagnostics methods implicitly focus on finding and avoiding small eigenvalues of $\mathbf{X}'\mathbf{X}$. The spectral decomposition of the sums of squares and cross products matrix (SSCP, Definition 1.44) is

$$\mathbf{X}'\mathbf{X} = \mathbf{V}\text{Dg}(\mathbf{d})\mathbf{V}' \tag{1.109}$$

If \mathbf{X} has full rank, then $(\mathbf{X}'\mathbf{X})^{-1} = \mathbf{V}[\text{Dg}(\mathbf{d})]^{-1}\mathbf{V}'$ exists and

$$\left|(\mathbf{X}'\mathbf{X})^{-1}\right| = (|(\mathbf{X}'\mathbf{X})|)^{-1} = \left(\prod_{k=1}^q d_k\right)^{-1} \tag{1.110}$$

As $\min(d_k) \downarrow 0$ the determinant of the inverse $\uparrow \infty$. If \mathbf{X} is less than full rank, then $(\mathbf{X}'\mathbf{X})^+ = \mathbf{V}[\text{Dg}(\mathbf{d})]^+\mathbf{V}'$ and \mathbf{d} includes some zeros. Equivalently $(\mathbf{X}'\mathbf{X})^+ = \mathbf{V}\text{Dg}(\mathbf{d}^{-1}, \mathbf{0})\mathbf{V}'$. The condition of the numerical properties of the matrix may be judged in terms of condition values, such as $\sqrt{\max(\mathbf{d})/\min(\mathbf{d})}$.

1.13 SOME FACTORS OF SYMMETRIC MATRICES

For a real and nonnegative number, such as $a = 25$, one can find its square root, \sqrt{a} , such that $(\sqrt{a})(\sqrt{a}) = a$. A diagonal matrix of nonnegative numbers $\mathbf{D} = \text{Dg}(\{d_1, \dots, d_p\})$ allows defining $\mathbf{D}^{1/2} = \text{Dg}(\{d_1^{1/2}, \dots, d_p^{1/2}\})$, with $\mathbf{D}^{1/2}\mathbf{D}^{1/2} = \mathbf{D}$. The concept extends to square and symmetric matrices in a variety of ways.

Definition 1.39 For $p \times p$ and symmetric \mathbf{A} , a $p_1 \times p$ factor with $p_1 \leq p$ is any \mathbf{F} such that $\mathbf{A} = \mathbf{F}\mathbf{F}'$.

The definition does not suffice to guarantee a unique factor, even with full rank \mathbf{A} . Depending on the side conditions desired, more than one factoring exists. In some analytic or computational settings any choice will serve, while other settings demand a particular choice. Full rank \mathbf{A} does always imply any factor is full rank, necessarily $p \times p$, and $\mathbf{A}^{-1} = \mathbf{F}^{-t}\mathbf{F}^{-1}$. More generally, $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{A})$, but it may not have the same dimensions when \mathbf{A} is not full rank.

Theorem 1.13 Any $p \times p$ and symmetric \mathbf{A} can be factored in three ways.

(a) The *square root*, or *Cholesky*, factor adds the side condition of lower triangular factor:

$$\mathbf{A} = \mathbf{L}\mathbf{L}', \tag{1.111}$$

with \mathbf{L} $p \times p$, lower triangular, with nonnegative diagonal elements. If $\text{rank}(\mathbf{A}) = p$, then \mathbf{L} is also full rank with strictly positive diagonal elements.

(b) The spectral decomposition, $\mathbf{A} = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})\mathbf{V}'$, implies a distinct factoring. If

$$\mathbf{F} = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})^{1/2}, \tag{1.112}$$

then $\mathbf{A} = \mathbf{F}\mathbf{F}'$, with $\text{Dg}(\boldsymbol{\lambda})^{1/2} = \text{Dg}(\{\lambda_j^{1/2}\})$ and $\mathbf{V}\mathbf{V}' = \mathbf{V}'\mathbf{V} = \mathbf{I}_p$ (without loss of generality). Like \mathbf{A} , the factor \mathbf{F} is $p \times p$ and $\text{rank}(\mathbf{F}) = \text{rank}(\mathbf{A})$.

(c) If $\text{rank}(\mathbf{A}) = p_1 \leq p$ and $\boldsymbol{\lambda}_1$ indicates the vector of p_1 nonzero eigenvalues, then without loss of generality we may choose $\mathbf{F} = [\mathbf{V}_+ \mathbf{V}_0]\text{Dg}(\boldsymbol{\lambda}_1, \mathbf{0})^{1/2}$, which is $p \times p$ and $\text{rank } p_1$. In contrast, the matrix

$$\mathbf{F}_1 = \mathbf{V}_1\text{Dg}(\boldsymbol{\lambda}_1)^{1/2} \tag{1.113}$$

is $p \times p_1$ and $\text{rank } p_1$. In turn

$$\begin{aligned} \mathbf{A} &= \mathbf{F}_1\mathbf{F}_1' \\ &= \mathbf{V}_1\text{Dg}(\boldsymbol{\lambda}_1)\mathbf{V}_1'. \end{aligned} \tag{1.114}$$

(d) Except in special cases, \mathbf{L} , \mathbf{F} and \mathbf{F}_1 are not symmetric. In some situations a symmetric factor may be preferred. If so, then choosing

$$\mathbf{F}_s = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})^{1/2}\mathbf{V}' \tag{1.115}$$

ensures $\mathbf{F}_s = \mathbf{F}_s'$ and $\mathbf{F}_s\mathbf{F}_s' = \mathbf{F}_s'\mathbf{F}_s = \mathbf{F}_s\mathbf{F}_s = \mathbf{A}$, with \mathbf{F}_s $p \times p$ and $\text{rank}(\mathbf{F}_s) = p_1$.

(e) An alternative symmetric factor is given by

$$\mathbf{F}_{s+} = \mathbf{V}_1\text{Dg}(\boldsymbol{\lambda}_1)^{1/2}\mathbf{V}_1', \tag{1.116}$$

with \mathbf{F}_{s+} $p \times p$ and $\text{rank}(\mathbf{F}_{s+}) = p_1$.

(f) If \mathbf{A} has any negative eigenvalues, complex variables will occur in the factor.

We use the theorem most often to factor a covariance matrix. The spectral factor becomes particularly convenient in proofs with less-than-full-rank covariance. The nice properties of diagonal and orthonormal matrices lead to simple expressions for useful generalized inverses, including $(\mathbf{F}_1)^+ = \text{Dg}(\boldsymbol{\lambda}_1)^{-1/2}\mathbf{V}_1'$. In contrast, for computational purposes the Cholesky factor and related QR decomposition should always be used.

1.14 SINGULAR VALUE DECOMPOSITION

Theorem 1.14 (a) Any $m \times n$ matrix \mathbf{A} has a *singular value decomposition*, $\mathbf{A} = \mathbf{U}\mathbf{S}\mathbf{V}'$, with orthonormal $m \times m$ \mathbf{U} of rank m , orthonormal $n \times n$ \mathbf{V} of rank n , and $s_{jk} = 0$ except for the main diagonal elements, $\{s_{jj}\}$, which are nonnegative. If $m \geq n \geq n_1 = \text{rank}(\mathbf{A})$, then \mathbf{s} has n_1 strictly positive values.

If $m \geq n \geq n_1$, then

$$\mathbf{A} = \mathbf{USV}' \quad (1.117)$$

$$\mathbf{U}'\mathbf{U} = \mathbf{UU}' = \mathbf{I}_m \quad (1.118)$$

$$\mathbf{V}'\mathbf{V} = \mathbf{VV}' = \mathbf{I}_n \quad (1.119)$$

$$\mathbf{S} = \begin{bmatrix} \text{Dg}(\mathbf{s}) \\ \mathbf{0}_{(m-n) \times n} \end{bmatrix}, \quad (1.120)$$

If $m < n$, then a similar construction applies to \mathbf{A}' .

(b) The vector \mathbf{s} contains the *singular values*, some or all of which may be zero. Necessarily $\text{Dg}(\mathbf{s})^2 = \text{Dg}(\{s_i^2\})$. Singular values equal the positive square roots of the eigenvalues of $n \times n \mathbf{A}'\mathbf{A}$, which coincide, except perhaps for some zeros, with the eigenvalues of $m \times m \mathbf{AA}'$.

(c) Without loss of generality, elements of \mathbf{s} may be assumed to be sorted from largest to smallest. If so, and $m \geq n \geq n_1 = \text{rank}(\mathbf{A})$, then $\mathbf{s}' = [\mathbf{s}'_1 \ \mathbf{0}'_{n-n_1}]$ with \mathbf{s}_1 the $n_1 \times 1$ vector of strictly positive elements. In turn $\mathbf{U} = [\mathbf{U}_1 \ \mathbf{U}_0]$ with \mathbf{U}_1 $m \times n_1$, $\mathbf{U}'_1\mathbf{U}_1 = \mathbf{I}_{n_1}$, $\mathbf{V} = [\mathbf{V}_1 \ \mathbf{V}_0]$ with \mathbf{V}_1 $n \times n_1$, $\mathbf{V}'_1\mathbf{V}_1 = \mathbf{I}_{n_1}$ and

$$\begin{aligned} \mathbf{A} &= [\mathbf{U}_1 \ \mathbf{U}_0] \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}'_1 \\ \mathbf{V}'_0 \end{bmatrix} \\ &= [\mathbf{U}_1 \text{Dg}(\mathbf{s}_1) \ \mathbf{0}_{m \times (n-n_1)}] \begin{bmatrix} \mathbf{V}'_1 \\ \mathbf{V}'_0 \end{bmatrix} \\ &= \mathbf{U}_1 \text{Dg}(\mathbf{s}_1) \mathbf{V}'_1 + \mathbf{0}_{m \times n} \\ &= \mathbf{U}_1 \text{Dg}(\mathbf{s}_1) \mathbf{V}'_1. \end{aligned} \quad (1.121)$$

(d) Without loss of generality, \mathbf{V} and \mathbf{U} may be chosen such that

$$\mathbf{A}'\mathbf{A} = \mathbf{VDg}(\mathbf{s})^2\mathbf{V}' = \mathbf{V}_1\text{Dg}(\mathbf{s}_1)^2\mathbf{V}'_1 \quad (1.122)$$

and

$$\mathbf{AA}' = \mathbf{UDg}(\mathbf{s}, \mathbf{0}_{m-n})^2\mathbf{U}' = \mathbf{U}_1\text{Dg}(\mathbf{s}_1)^2\mathbf{U}'_1. \quad (1.123)$$

In practice, the particular choices for \mathbf{U} and \mathbf{V} are intertwined due to the fact that an eigenvector remains an eigenvector if it is multiplied by -1 . In order to account for the signs, having chosen \mathbf{V}_1 requires choosing $\mathbf{U}_1 = \mathbf{AV}_1\text{Dg}(\mathbf{s}_1)^{-1}$. Alternately, having chosen \mathbf{U}_1 requires choosing $\mathbf{V}'_1 = \text{Dg}(\mathbf{s}_1)^{-1}\mathbf{U}'_1\mathbf{A}$. If needed, the definitions of \mathbf{U} and \mathbf{V} may be completed by adding the eigenvectors which correspond to zero eigenvalues of \mathbf{AA}' and $\mathbf{A}'\mathbf{A}$.

Lemma 1.28 For any $m \times n$ matrix \mathbf{A} , the $n \times m$ Moore-Penrose (four-condition, unique) generalized inverse is

$$\begin{aligned}
 \mathbf{A}^+ &= \mathbf{V}\mathbf{S}^+\mathbf{U}' \\
 &= \mathbf{V} \left[\text{Dg}(\mathbf{s})^+ \quad \mathbf{0}_{n \times (m-n)} \right] \mathbf{U}' \\
 &= \mathbf{V}_1 [\text{Dg}(\mathbf{s}_1)]^{-1} \mathbf{U}'_1.
 \end{aligned}
 \tag{1.124}$$

Of course $[\text{Dg}(\mathbf{s})]^+ = \text{Dg}(\{s_1^{-1}, \dots, s_r^{-1}, 0, \dots, 0\})$.

The singular value and spectral decomposition can be chosen to coincide for square, symmetric, positive definite and semidefinite matrices (such as inner and outer product matrices).

1.15 PROJECTIONS AND OTHER FUNCTIONS OF A DESIGN MATRIX

Design matrices play a central role in the study of linear models. A design matrix, \mathbf{X} , for N observations and q variables is of dimension $N \times q$ with $N \geq q$ and $\text{rank}(\mathbf{X}) = r \leq q$. Many times throughout the book we shall take advantage of various properties of the singular value decomposition (SVD) of \mathbf{X} , as well as the closely related spectral decompositions of $\mathbf{X}'\mathbf{X}$ and $\mathbf{X}\mathbf{X}'$. The following lemma summarizes the properties in an integrated presentation.

Lemma 1.29 (a) For $N \times q$ matrix \mathbf{X} , with $N \geq q$ and $\text{rank}(\mathbf{X}) = r \leq q$, the SVD gives

$$\begin{aligned}
 \mathbf{X} &= \mathbf{L} \begin{bmatrix} \text{Dg}(\mathbf{s}) \\ \mathbf{0} \end{bmatrix} \mathbf{R}' \\
 &= [\mathbf{L}_1 \quad \mathbf{L}_0] \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{R}'_1 \\ \mathbf{R}'_0 \end{bmatrix} \\
 &= \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \mathbf{R}'_1.
 \end{aligned}
 \tag{1.125}$$

The first product has dimensions $(N \times N)(N \times q)(q \times q)$, and the last product has dimensions $(N \times r)(r \times r)(r \times q)$. Without loss of generality, singular values are sorted from largest to smallest, $\mathbf{s} = [\mathbf{s}'_1 \quad \mathbf{0}_{n-n_1}]$, \mathbf{s}_1 is $n_1 \times 1$ with all strictly positive values. Also $\mathbf{R} = [\mathbf{R}_1 \quad \mathbf{R}_0]$, $\mathbf{R}'\mathbf{R} = \mathbf{R}\mathbf{R}' = \mathbf{I}_q$ with \mathbf{R}_1 $q \times n_1$. Similarly $\mathbf{L} = [\mathbf{L}_1 \quad \mathbf{L}_0]$ with \mathbf{L}_1 $n \times n_1$ and $\mathbf{L}'\mathbf{L} = \mathbf{L}\mathbf{L}' = \mathbf{I}_N$.

(b) Here \mathbf{R}_1 and \mathbf{R}_0 are sets of orthonormal eigenvectors (unique, up to reflections) corresponding to nonzero (\mathbf{R}_1) and zero (\mathbf{R}_0) eigenvalues of $\mathbf{X}'\mathbf{X}$. Specifically, without loss of generality, SVD properties correspond to assuming

$$\begin{aligned}
 \mathbf{X}'\mathbf{X} &= \mathbf{R}\text{Dg}(\mathbf{s})^2\mathbf{R}' \\
 &= \mathbf{R}_1\text{Dg}(\mathbf{s}_1)^2\mathbf{R}'_1.
 \end{aligned}
 \tag{1.126}$$

(c) Similarly, \mathbf{L}_1 and \mathbf{L}_0 are (nonunique) sets of orthonormal eigenvectors corresponding to nonzero (\mathbf{L}_1) and zero (\mathbf{L}_0) eigenvalues of $\mathbf{X}\mathbf{X}'$. Specifically, without loss of generality, SVD properties correspond to assuming

$$\begin{aligned} \mathbf{X}\mathbf{X}' &= \mathbf{L} \begin{bmatrix} \text{Dg}(\mathbf{s}_1, \mathbf{0})^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{L}' \\ &= \mathbf{L}_1 \text{Dg}(\mathbf{s}_1)^2 \mathbf{L}'_1. \end{aligned} \quad (1.127)$$

(d) Given \mathbf{R}_1 and \mathbf{s}_1 , necessarily

$$\mathbf{L}_1 = \mathbf{X}\mathbf{R}_1 \text{Dg}(\mathbf{s}_1)^{-1}. \quad (1.128)$$

Alternately, given \mathbf{R}_1 and \mathbf{s}_1 , necessarily

$$\mathbf{R}_1 = \text{Dg}(\mathbf{s}_1)^{-1} \mathbf{L}'_1 \mathbf{X}. \quad (1.129)$$

However, only one of \mathbf{L}_1 and \mathbf{R}_1 can be chosen freely, due to the need to account for the sign ambiguity of eigenvectors (as detailed in Theorem 1.12).

(e) Furthermore

$$\begin{aligned} (\mathbf{X}'\mathbf{X})^+ &= \mathbf{R} [\text{Dg}(\mathbf{s})^2]^+ \mathbf{R}' \\ &= \mathbf{R}_1 \text{Dg}(\mathbf{s}_1)^{-2} \mathbf{R}'_1. \end{aligned} \quad (1.130)$$

Proof. The results follow from the SVD definition, based on properties of orthonormal matrices, inverses, and eigenanalysis of symmetric matrices.

Many functions of design matrices, especially ones involving $(\mathbf{X}'\mathbf{X})^-$, must be understood to develop the theory of linear models. The value of the results lies in guaranteeing uniqueness, symmetry, and the ability to simplify expressions in estimators. Even though $(\mathbf{X}'\mathbf{X})$ is always symmetric, $(\mathbf{X}'\mathbf{X})^-$ need not be symmetric. Furthermore a less-than-full-rank \mathbf{X} leads to infinitely many choices for $(\mathbf{X}'\mathbf{X})^-$, while $(\mathbf{X}'\mathbf{X})^+$ is unique. The theorem gives many equalities useful for simplification and many invariance properties guaranteeing uniqueness. The corollary summarizes many special properties associated with the matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$, which plays a key role in linear models estimation and distribution theory.

Definition 1.40 (a) The matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$ is a *projection matrix*.

(b) Matrix $\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{Y}$ is the *projection of Y* into the space spanned by the columns of \mathbf{X} .

Theorem 1.15 For any matrix \mathbf{X} , the following all hold.

- (a) $[(\mathbf{X}'\mathbf{X})^-]'$ is also a generalized inverse of $(\mathbf{X}'\mathbf{X})$.
- (b) $\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X} = \mathbf{X}$; hence $(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$ is a generalized inverse of \mathbf{X} .
- (c) $\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$ is invariant to $(\mathbf{X}'\mathbf{X})^-$.
- (d) $\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$ is always symmetric, even if $(\mathbf{X}'\mathbf{X})^-$ is not symmetric.
- (e) $\mathbf{X}[(\mathbf{X}'\mathbf{X})^-]'\mathbf{X}'\mathbf{X} = \mathbf{X}$.
- (f) $\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}' = \mathbf{X}'$.
- (g) $\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^-]'\mathbf{X}' = \mathbf{X}'$.

- (h) $\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-}]'\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$.
- (i) $\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-}]'\mathbf{X}'$ is always symmetric even if $[(\mathbf{X}'\mathbf{X})^{-}]'$ is not symmetric.
- (j) $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{+}\mathbf{X}'$.
- (k) All of the matrices mentioned in the theorem have rank $r = \text{rank}(\mathbf{X})$.

Proof. Searle (1971) stated and proved parts (a)–(d) as his Theorem 7 and parts (e)–(i) as a corollary. Part (j) follows from uniqueness of the four-condition inverse. Part (k) follows from four- and one- condition inverse properties, as well as Lemma 1.29. \square

Corollary 1.15 For any $N \times q$ matrix \mathbf{X} with $\text{rank}(\mathbf{X}) = r \leq q \leq N$ the projection matrix, $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$, is (a) unique, (b) symmetric, (c) idempotent and (d) rank r , with (e) r eigenvalues of one and $N - r$ of zero. Furthermore, $\mathbf{I} - \mathbf{H}$ is (f) unique, (g) symmetric, (h) idempotent and (i) rank $N - r$, with (j) $N - r$ eigenvalues of one and r of zero. (k) Also, $\mathbf{H}(\mathbf{I} - \mathbf{H}) = \mathbf{0}$. (l) With exactly the notation of Lemma 1.29,

$$\begin{aligned} \mathbf{H} &= \mathbf{L} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{N-r} \end{bmatrix} \mathbf{L}' \\ &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{N-r} \end{bmatrix} \begin{bmatrix} \mathbf{L}'_1 \\ \mathbf{L}'_0 \end{bmatrix} \\ &= \mathbf{L}_1 \mathbf{L}'_1 \end{aligned} \tag{1.131}$$

and

$$\begin{aligned} \mathbf{I} - \mathbf{H} &= \mathbf{L} \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{N-r} \end{bmatrix} \mathbf{L}' \\ &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{N-r} \end{bmatrix} \begin{bmatrix} \mathbf{L}'_1 \\ \mathbf{L}'_0 \end{bmatrix} \\ &= \mathbf{L}_0 \mathbf{L}'_0. \end{aligned} \tag{1.132}$$

Here $\mathbf{L}'\mathbf{L} = \mathbf{L}\mathbf{L}' = \mathbf{I}_N$, while $\mathbf{L}'_1\mathbf{L}_1 = \mathbf{I}_r$ and $\mathbf{L}'_0\mathbf{L}_0 = \mathbf{I}_{N-r}$. (m) Finally, \mathbf{L} contains orthonormal eigenvectors corresponding to nonzero (\mathbf{L}_1) and zero (\mathbf{L}_0) eigenvalues of $\mathbf{X}\mathbf{X}'$:

$$\begin{aligned} \mathbf{X}\mathbf{X}' &= \mathbf{L} \begin{bmatrix} \text{Dg}(\mathbf{s}_1, \mathbf{0})^2 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \mathbf{L}' \\ &= \mathbf{L}_1 \text{Dg}(\mathbf{s}_1)^2 \mathbf{L}'_1. \end{aligned} \tag{1.133}$$

Proof. Results follow from combining forms in the theorem and Lemma 1.29. Writing $\mathbf{H}^2 = (\mathbf{L}_1 \mathbf{L}'_1)(\mathbf{L}_1 \mathbf{L}'_1) = \mathbf{L}_1 \mathbf{L}'_1$ demonstrates that \mathbf{H} is idempotent. \square

The matrix $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{+}\mathbf{X}'$ earned the nickname “hat matrix” in the linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ because $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$. Uniqueness of \mathbf{H}

implies uniqueness of many functions of H , including \hat{y} , as well as $\hat{e} = (I - H)y$ and $\hat{\sigma}^2 = \hat{e}'\hat{e}/(N - r)$.

Lemma 1.30 If A ($n \times n$) is a constant matrix and F' ($n \times n$) is nonsingular with $B = FF'$, then $F'AF$ is idempotent $\Leftrightarrow AB$ is idempotent.

Proof. $FF' = B$ implies $AFF' = AB$. Premultiplying by F' and postmultiplying by F^{-t} gives $F'AF = F'ABF^{-t}$. Hence $(F'AF)^2 = F'ABABF^{-t}$. If $(F'AF)^2 = F'AF$ then $(AB)^2 = AB$, which implies $(F'AF)^2 = (F'AF)$. \square

Definition 1.41 Matrix P ($N \times N$) is a *permutation matrix* if it is obtained by permuting the rows (only) of I_N .

Given the definition, premultiplying by a conforming permutation matrix exchanges rows. Postmultiplying by a conforming transposed permutation matrix exchanges columns. A permutation matrix is always orthonormal and full rank, which implies $P^{-1} = P'$. A permutation matrix may always be found which, when multiplied times a conforming matrix, permutes the rows (or columns) to any new order desired.

1.16 SPECIAL PROPERTIES OF PATTERNED MATRICES

A partitioned matrix may be thought of as a supermatrix, a matrix containing matrices. Most importantly, if the partitions conform for addition and matrix multiplication then the results can be expressed in terms of the partitions, without considering particular elements. *It is crucial to verify conformation of each pair of partitions as well as the total matrices.* When the partitions do conform, the rules of matrix multiplication apply to create what might be thought of as “super” or “meta” multiplication or addition. Examples include the following (when the matrices conform for the operations):

$$A[B_1 \ B_2] = [AB_1 \ AB_2] \tag{1.134}$$

$$[B_1 \ B_2] \begin{bmatrix} C_1 \\ C_2 \end{bmatrix} = B_1C_1 + B_2C_2 \tag{1.135}$$

$$[B_1 \ B_2][B_1' \ B_2'] = B_1B_1' + B_2B_2' \tag{1.136}$$

$$[B_1 \ B_2]'[B_1 \ B_2] = \begin{bmatrix} B_1'B_1 & B_1'B_2 \\ B_2'B_1 & B_2'B_2 \end{bmatrix}. \tag{1.137}$$

Not only multiplication and addition, but also determinants and inverses often can be expressed compactly and conveniently in terms of partitions.

Theorem 1.16 If $p \times p$ A is positive definite and symmetric, $B = A^{-1}$, and A and B are partitioned with corresponding submatrices of the same sizes,

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix} \tag{1.138}$$

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}_{21} & \mathbf{B}_{22} \end{bmatrix}, \tag{1.139}$$

then (given the indicated operations are valid)

$$\mathbf{B}_{11} = (\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21})^{-1} = \mathbf{A}_{11}^{-1} + \mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{B}_{22}\mathbf{A}_{21}\mathbf{A}_{11}^{-1} \tag{1.140}$$

$$\mathbf{B}_{12} = -\mathbf{A}_{11}^{-1}\mathbf{A}_{12}\mathbf{B}_{22} \tag{1.141}$$

$$\mathbf{B}_{21} = -\mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{B}_{11} \tag{1.142}$$

$$\mathbf{B}_{22} = (\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12})^{-1} = \mathbf{A}_{22}^{-1} + \mathbf{A}_{22}^{-1}\mathbf{A}_{21}\mathbf{B}_{11}\mathbf{A}_{12}\mathbf{A}_{22}^{-1}. \tag{1.143}$$

The subscript pattern gives a mnemonic device for remembering the formulas.

Theorem 1.17 For any $p \times p$ \mathbf{A} partitioned as

$$\mathbf{A} = \begin{bmatrix} \mathbf{A}_{11} & \mathbf{A}_{12} \\ \mathbf{A}_{21} & \mathbf{A}_{22} \end{bmatrix}, \tag{1.144}$$

(a) if either $\mathbf{A}_{12} = \mathbf{0}$ or $\mathbf{A}_{21} = \mathbf{0}$, then $|\mathbf{A}| = |\mathbf{A}_{11}||\mathbf{A}_{22}|$.

(b) For any (conforming) $\{\mathbf{A}_{11}, \mathbf{A}_{12}, \mathbf{A}_{21}\}$, if \mathbf{A}_{22} is full rank, then $|\mathbf{A}| = |\mathbf{A}_{22}||\mathbf{A}_{11} - \mathbf{A}_{12}\mathbf{A}_{22}^{-1}\mathbf{A}_{21}|$.

(c) For any \mathbf{A}_{12} or \mathbf{A}_{21} , if \mathbf{A}_{11} is full rank, then $|\mathbf{A}| = |\mathbf{A}_{11}||\mathbf{A}_{22} - \mathbf{A}_{21}\mathbf{A}_{11}^{-1}\mathbf{A}_{12}|$.

Theorem 1.18 (a) If \mathbf{A} and \mathbf{B} are the same size then

$$(\mathbf{A} + \mathbf{B}) \otimes \mathbf{C} = (\mathbf{A} \otimes \mathbf{C}) + (\mathbf{B} \otimes \mathbf{C}) \tag{1.145}$$

and

$$\mathbf{C} \otimes (\mathbf{A} + \mathbf{B}) = (\mathbf{C} \otimes \mathbf{A}) + (\mathbf{C} \otimes \mathbf{B}). \tag{1.146}$$

(b) If $\mathbf{A}, \mathbf{B}, \mathbf{C}$, and \mathbf{D} are $m \times h, p \times k, h \times n$ and $k \times q$, respectively, then

$$(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC}) \otimes (\mathbf{BD}). \tag{1.147}$$

1.17 FUNCTION OPTIMIZATION AND MATRIX DERIVATIVES

Deriving properties of linear models often leads to the need to maximize or minimize a smooth function. In addition, side conditions may be desired which impose constraints on the optimization. The most convenient approach usually involves creating a system of equations to be solved by determining derivatives of the function. In turn, the introduction of Lagrangian multipliers often satisfies the need to impose side conditions.

Given the focus of the book, we assume the reader has substantial knowledge of the theory and practical use of derivatives for the analysis of scalar functions. In the few instances in the book when we need to know specific vector and matrix forms of derivatives we present the required results without comment. A defensible description of even the basic rules of matrix derivatives requires attention to concepts in real analysis rather far from the topic at hand. We recommend the reader seeking the motivation for the derivatives used here consult the excellent book by Magnus and Neudecker (1988) for a definitive treatment. Schott (2005) provided a more brief but carefully constructed introduction.

1.18 STATISTICAL NOTATION INVOLVING MATRICES

A common notation for the study of probability uses Greek letters for parameters, uppercase Roman letters for random variables, and lowercase Roman letters for realizations of random variables (particular sample values). The convention conflicts with matrix notation distinctions between scalars, vectors, and matrices. The importance of matrix expressions throughout the book means that matrix notation must dominate. Consequently the reader must often distinguish fixed from random, known from unknown, observed from unobserved, and observable from unobservable via the context of the discussion. When in doubt about a particular item, simply search backwards in the text to discover where the variable was introduced.

Expressions involving random variables use the portion of conventional notation that does not conflict with the matrix notation introduced in Section 1.1. Roman letters towards the beginning of the alphabet, such as $\{c, e, C\}$, will usually represent constants (either known or unspecified). In contrast, Roman letters towards the end of the alphabet, such as $\{y, \mathbf{y}, Y\}$, will usually represent random variables (which may take on infinitely many values). In turn, $\{y_*, \mathbf{y}_*, Y_*\}$ will represent an arbitrary possible value, such as a variable of integration, while $\{y_0, \mathbf{y}_0, Y_0\}$ will represent a single but unspecified particular value. For two jointly absolutely continuous random variables, $\{x, y\}$ the rules lead to the expression $E(x|y = y_0) = \int x_* f_{x,y}(x_*, y_0) dx_*$. Greek letters, such as $\{\beta, \boldsymbol{\beta}, B\}$, will represent parameters (fixed and unknown population properties). Corresponding (random) estimators will be indicated $\{\hat{\beta}, \hat{\boldsymbol{\beta}}, \hat{B}\}$ or $\{\tilde{\beta}, \tilde{\boldsymbol{\beta}}, \tilde{B}\}$.

1.19 STATISTICAL FORMULAS

With $N \gg p$, matrix $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \cdots \ \mathbf{y}_p]$ contains a collection of scores on p variables. Rows are independent sampling units and columns are variables. Although discussed in more detail in Chapter 7, the concepts of mean and variance for a vector are useful here.

Definition 1.42 (a) For $\mathbf{y} = \{y_j\}$ a vector of random variables with a well-defined joint distribution $E(\mathbf{y}) = \{E(y_j)\}$ defines the *population mean*, when $E(y_j)$ exists $\forall j$.

(b) Similarly, $E(\mathbf{Y}) = \{E(y_{jk})\}$ when it exists.

(c) For $\mathbf{y} = \{y_j\}$ a vector of random variables with a well-defined joint distribution and $E(\mathbf{y}) = \boldsymbol{\mu}$, $\mathcal{V}(\mathbf{y}) = \{E[(y_j - \mu_j)(y_k - \mu_k)]\}$ defines the *population covariance matrix*, when $E[(y_j - \mu_j)(y_k - \mu_k)]$ exists $\forall (j, k)$.

The concept of $\mathcal{V}(\mathbf{Y})$ is not well defined. It is customary and fully satisfactory to consider either $\mathcal{V}[\text{vec}(\mathbf{Y})]$ or $\mathcal{V}[\text{vec}(\mathbf{Y}^*)]$, which are well defined. They differ only by a permutation of rows, which can be achieved by multiplication with a permutation matrix. The following definition and lemma are repeated and discussed in Chapter 7. They are presented here to allow a precise description of principal components analysis of random data.

Definition 1.43 (a) A random vector \mathbf{y} ($p \times 1$) with finite second moments has an associated *covariance matrix*

$$\mathcal{V}(\mathbf{y}) = \{E(y_j y_k)\} - \{E(y_j)E(y_k)\} = \boldsymbol{\Sigma}. \quad (1.148)$$

(b) If all $\sigma_{jj} = \langle \boldsymbol{\Sigma} \rangle_{jj}$ are such that $0 < \sigma_{jj} < \infty$, then \mathbf{y} has *correlation matrix*

$$\begin{aligned} \mathbf{P} &= \text{Dg}(\{\sigma_{11}, \dots, \sigma_{pp}\})^{-1/2} \boldsymbol{\Sigma} \text{Dg}(\{\sigma_{11}, \dots, \sigma_{pp}\})^{-1/2} \\ &= \{\rho_{jk}\}. \end{aligned} \quad (1.149)$$

Lemma 1.31 Finite second moments for a random vector \mathbf{y} ($p \times 1$) guarantee the existence of finite covariance and correlation matrices for the population.

(a) The population covariance matrix can be expressed as

$$\begin{aligned} \mathcal{V}(\mathbf{y}) &= \{E(y_j y_k)\} - \{E(y_j)E(y_k)\} \\ &= \{E(y_j y_k)\} - \{\mu_j \mu_k\} \\ &= E(\mathbf{y} \mathbf{y}') - \boldsymbol{\mu} \boldsymbol{\mu}' \\ &= E(\mathbf{y} \mathbf{y}') - E(\mathbf{y})[E(\mathbf{y})]' \\ &= \boldsymbol{\Sigma}. \end{aligned} \quad (1.150)$$

It is symmetric and nonnegative definite (positive definite or positive semidefinite). The covariance matrix contains centered, average cross products, with diagonal element σ_{jj} the variance of y_j .

(b) If $\sigma_{jj} > 0 \forall j$, then the population correlation matrix \mathbf{P} is well defined and nonnegative definite, with $\rho_{jk} = \sigma_{jk} / \sqrt{\sigma_{jj} \sigma_{kk}}$. The correlation matrix contains centered and scaled average cross products.

Definition 1.44 (a) A set of N observations of a $p \times 1$ random vector arranged in an $N \times p$ matrix, with each observation forming a row, provides a *data matrix*, such as \mathbf{Y} .

(b) With $N \times 1$ vector $\mathbf{1}_N = [1 \ 1 \ \cdots \ 1]'$, the $p \times 1$ *sample mean vector* is

$$\bar{\mathbf{y}} = \mathbf{Y}'\mathbf{1}_N/N. \quad (1.151)$$

(c) The $p \times p$ *sample SSCP matrix* $\mathbf{Y}'\mathbf{Y}$ contains sums of squares and cross products.

(d) The $p \times p$ *sample covariance matrix* is

$$\begin{aligned} \mathbf{S} &= (\mathbf{Y} - \mathbf{1}_N\bar{\mathbf{y}})'(\mathbf{Y} - \mathbf{1}_N\bar{\mathbf{y}})/N \\ &= \mathbf{Y}'[\mathbf{I}_N - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}']\mathbf{Y}/N \\ &= \mathbf{Y}'\mathbf{Y}/N - \bar{\mathbf{y}}\bar{\mathbf{y}}'. \end{aligned} \quad (1.152)$$

(e) If $s_{jj} > 0$, then the $p \times p$ *sample correlation matrix* is

$$\mathbf{R} = \text{Dg}(\{s_{11}, \dots, s_{pp}\})^{-1/2} \mathbf{S} \text{Dg}(\{s_{11}, \dots, s_{pp}\})^{-1/2}, \quad (1.153)$$

with $r_{jk} = s_{jk}/\sqrt{s_{jj}s_{kk}}$.

Lemma 1.32 (a) If data matrix \mathbf{Y} has independent rows with $E[\text{row}_i(\mathbf{Y})] \equiv \boldsymbol{\mu}'$ and $\mathcal{V}[\text{row}_i(\mathbf{Y})] = \boldsymbol{\Sigma}$, then

$$\widehat{\boldsymbol{\Sigma}} = \mathbf{S}N/(N-1) \quad (1.154)$$

and $E(\widehat{\boldsymbol{\Sigma}}) = \boldsymbol{\Sigma}$. More generally, in fitting multivariate linear models, $\widehat{\boldsymbol{\Sigma}} = \mathbf{S}N/(N-r)$, in which r equals the rank of the design matrix, provides an unbiased estimator when $N > r$.

(b) If the sample correlation matrix estimates the population matrix, $\mathbf{R} = \widehat{\mathbf{P}}$, then $E(\widehat{\mathbf{P}}) \neq \mathbf{P}$, except when $\mathbf{P} = \mathbf{I}_p$ (which implies $\rho_{jk} = 0$ if $j \neq k$). If $\mathbf{P} \neq \mathbf{I}_p$, then no constant c can be found to make $c\widehat{\mathbf{P}}$ unbiased.

If r_{-jj} indicates diagonal element j of \mathbf{R}^{-1} , then $R_j^2 = (r_{-jj} - 1)/r_{-jj}$ is the squared multiple correlation between variable j and the remaining $p-1$ variables.

Partitioning the variables into two sets, with $\mathbf{Y} = [\mathbf{Y}_1 \ \mathbf{Y}_2]$, gives p_1 variable in \mathbf{Y}_1 and p_2 variables in \mathbf{Y}_2 . A corresponding partitioning of the covariance matrix gives

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix}. \quad (1.155)$$

The sample covariance for $\mathbf{Y}_1|\mathbf{Y}_2$ equals

$$\mathbf{S}_{\mathbf{Y}_1|\mathbf{Y}_2} = \mathbf{S}_{11} - \mathbf{S}_{12}\mathbf{S}_{22}^{-1}\mathbf{S}_{21}, \quad (1.156)$$

with corresponding unbiased estimator $\mathbf{S}_{\mathbf{Y}_1|\mathbf{Y}_2}N/(N-p_2-1)$. If $\mathbf{D}_{\mathbf{Y}_1|\mathbf{Y}_2}$ indicates

the diagonal matrix containing the diagonal elements of $\mathbf{S}_{Y_1|Y_2}$, the corresponding correlation matrix equals

$$\mathbf{R}_{Y_1|Y_2} = \mathbf{D}_{Y_1|Y_2}^{-1/2} \mathbf{S}_{Y_1|Y_2} \mathbf{D}_{Y_1|Y_2}^{-1/2}. \tag{1.157}$$

Elements of $\mathbf{R}_{Y_1|Y_2}$ equal full partial correlations among the variables in \mathbf{Y}_1 because any two variables in \mathbf{Y}_1 have both been adjusted for variables in \mathbf{Y}_2 . Muller and Fetterman (2002, Chapter 6) gave a brief overview of various sorts of partial correlations in the context of univariate linear models regression.

1.20 PRINCIPAL COMPONENTS

Given a set of observations on a group of variables, it may be helpful to describe a simple model for the structure of the corresponding cross products, covariance, or correlation matrix. Three distinct applications may motivate the process: (1) analytic decomposition of population variables in a proof, (2) numerical analysis of observed values on a set of variables, and (3) data analysis for exploratory or confirmatory purposes. In the present section and the remainder of the book, we focus on the first application. Muller and Fetterman (2002, Chapter 8) provided a brief introduction to using principal components analysis for regression diagnostics. Timm (2002, Chapter 8) detailed the basic methods of principal components in the context of factor analysis models. Both texts also include useful additional references. Jackson (1991), and Jolliffe (2002) provided entire books about component analysis, while Basilevsky (1994) discussed component analysis within the more general context of factor analysis.

We strongly prefer factor analysis methods and related covariance structure models over principal components analysis for building and evaluating covariance models. One important reason arises from the concept of robustness to overfitting. A principal components analysis (PCA) model defines a special case of a factor analysis (FA) model. If the PCA model holds but the data analyst uses the FA model, then no harm should result. An adequate sample and analysis strategy should lead to reducing the model appropriately. In contrast, fitting a PCA model when the FA model holds always leads to bias and an invalid model, no matter how large the sample size. A parallel conclusion holds in seeking the best model for the response mean in a multiple regression. Overfitting only costs a bit of sensitivity, while underfitting creates bias and invalid models. Widaman (2004; also 1993) provided an extensive discussion of the question.

The discussion of principal components analysis will be cast primarily in terms of decomposing a population covariance matrix Σ . The underlying data are $N \times p$ $\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \cdots \ \mathbf{y}_p]$, with $N \gg p$, a collection of scores on p variables. Here rows of \mathbf{Y} are independent, with common covariance $\mathcal{V}[\text{row}_i(\mathbf{Y})] = \Sigma$. Most often, analysis involves \mathbf{Y} or $\mathbf{Y}_d = [\mathbf{I}_N - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}']\mathbf{Y} = \mathbf{Y} - \mathbf{1}\bar{\mathbf{y}}' = \{y_{ij} - \bar{y}_j\}$, the mean-centered transformation of \mathbf{Y} . With the impact of the means (first moments)

suppressed, the focus lies on modeling the second moments, which are variances and covariances.

Definition 1.45 (a) The first *principal component* equals the linear combination of a set of variables which has maximum variance, among all such combinations with unit-length coefficient vectors.
 (b) Such a linear combination of a set of variables defines a new variable, which will be called a *variate*.

Any covariance matrix has many special properties. A population or sample covariance matrix is always symmetric and can be expressed as an inner product. It always has a spectral decomposition with only positive or zero eigenvalues, which are the variances of the principal component variables. Each corresponding eigenvector holds a set of regression weights for the original variables which define the principal variables. The principal component variables are uncorrelated, with successively maximum variances (the eigenvalues).

The principal components have corresponding principal component scores, $\mathbf{Y}_c = [\mathbf{y}_{c,1} \ \mathbf{y}_{c,2} \ \cdots \ \mathbf{y}_{c,p}]$, an $N \times p$ matrix. The scores have many special properties. (1) Each column of \mathbf{Y}_c equals a linear combination of \mathbf{Y}_d : $\mathbf{Y}_c = \mathbf{Y}_d \mathbf{T}$. (2) $\mathcal{V}[\text{row}_i(\mathbf{Y}_c)] = \text{Dg}(\boldsymbol{\lambda})$. (3) The first set of N scores, $\mathbf{y}_{c,1}$, has maximum variance among all linear combinations (subject to the unit-length constraint $\mathbf{v}'_1 \mathbf{v}_1 = 1$), the second set of scores, $\mathbf{y}_{c,2}$, has maximum variance among all linear combinations given $\mathbf{y}_{c,1}$, etc. (4) The weight matrix $\mathbf{T} = [\mathbf{v}_1 \ \mathbf{v}_2 \ \cdots \ \mathbf{v}_p]$ must be the (orthonormal and full-rank) eigenvectors of the covariance matrix, $\boldsymbol{\Sigma} = \mathbf{T} \text{Dg}(\boldsymbol{\lambda}) \mathbf{T}'$ with $\mathbf{T}' \mathbf{T} = \mathbf{T} \mathbf{T}' = \mathbf{I}_p$. (5) For $k \neq k'$, the covariance and correlation between $\mathbf{y}_{c,k}$ and $\mathbf{y}_{c,k'}$ are zero.

Spectral decomposition of a sample estimator gives $\widehat{\boldsymbol{\Sigma}} = \widehat{\mathbf{T}} \text{Dg}(\widehat{\boldsymbol{\lambda}}) \widehat{\mathbf{T}}'$. For the particular set of sample values in hand, \mathbf{Y}_0 , which corresponds to $\widehat{\boldsymbol{\Sigma}}_0 = \widehat{\mathbf{T}}_0 \text{Dg}(\widehat{\boldsymbol{\lambda}}_0) \widehat{\mathbf{T}}_0'$, sample estimates attain optimal numerical properties in parallel to population properties. The component score estimator is $\widehat{\mathbf{Y}}_c = \mathbf{Y}_d \widehat{\mathbf{T}}$, with corresponding estimate $\widehat{\mathbf{Y}}_{0c} = \mathbf{Y}_{0d} \widehat{\mathbf{T}}_0$. Each $(p \times 1)$ column of $\widehat{\mathbf{T}}_0$ serves as the estimate of a set of p regression coefficients.

Here \mathbf{Y}_{0c} ($N \times p$) is a matrix with each column a set of component score estimates. Also $\bar{\mathbf{y}}_{0d} = \mathbf{0}$ implies $\bar{\mathbf{y}}_{0c} = \widehat{\mathbf{T}}_0 \bar{\mathbf{y}}_{0d} = \mathbf{0}$. The covariance matrix observed among the component score estimates is

$$\begin{aligned} \mathbf{S}_{0c} &= (\mathbf{Y}'_{0c} \mathbf{Y}_{0c} / N - \bar{\mathbf{y}}_{0c} \bar{\mathbf{y}}'_{0c}) N / (N - 1) \\ &= \mathbf{Y}'_{0c} \mathbf{Y}_{0c} / (N - 1) \\ &= \widehat{\mathbf{T}}_0' \mathbf{Y}'_{0d} \mathbf{Y}_{0d} \widehat{\mathbf{T}}_0 / (N - 1) \\ &= \widehat{\mathbf{T}}_0' \widehat{\boldsymbol{\Sigma}}_0 \widehat{\mathbf{T}}_0 \\ &= \text{Dg}(\widehat{\boldsymbol{\lambda}}_0). \end{aligned} \tag{1.158}$$

Consequently the estimated component scores $\mathbf{Y}_{0c} = \mathbf{Y}_{0d} \widehat{\mathbf{T}}_0$ are uncorrelated.

Also, the variance of the first component is $\hat{\lambda}_{0,1}$, the largest eigenvalue of $\hat{\Sigma}_0$, the variance of the second component is $\hat{\lambda}_{0,2}$, the second largest eigenvalue of $\hat{\Sigma}_0$, etc. The $p \times p$ factor (matrix) for $\hat{\Sigma}$, based on a spectral decomposition, equals $\hat{\Phi}_0 = \hat{\Upsilon}_0 \text{Dg}(\hat{\lambda}_0)^{1/2}$, with $\hat{\Sigma}_0 = \hat{\Phi}_0 \hat{\Phi}'_0$. The elements of $\hat{\Phi}_0 = \{\hat{\phi}_{0,jk}\}$ do not equal covariances between variable i and component j . The factor matrix is distinct from the eigenvectors (coefficients) and from the component scores.

1.21 SPECIAL COVARIANCE PATTERNS

Any $p \times p$ symmetric and nonnegative definite matrix provides a valid covariance matrix with up to $p(p + 1)$ distinct elements. Some special sampling schemes generate patterned covariance matrices, with elements expressible as a function of a small number of parameters. In the simplest case, complete independence (and finite second moments) gives $\Sigma = \sigma^2 I_p$. The study of time series leads to considering models such as autoregressive and moving average covariance patterns (Box, Jenkins, and Reinsel, 1994). The study of spatial statistics has also led to the development of a large range of covariance models (Cressie, 1991). The following definitions will be used throughout the book.

Definition 1.46 A set of jointly Gaussian variables with $\Sigma = \sigma^2 I_p$ have a *spherical* distribution in that equal-probability regions centered at the mean vector are circles for $p = 2$, spheres for $p = 3$, and hyperspheres for $p > 3$. In some contexts, *sphericity* is present when the weaker condition $\Sigma = V(\sigma^2 I_p)V'$ holds.

Definition 1.47 Any square, symmetric and nonnegative definite matrix with finite elements describes an *unstructured covariance matrix*.

Definition 1.48 A $p \times p$ compound symmetric covariance matrix may be expressed as $\Sigma = \sigma^2 [\mathbf{1}_p \mathbf{1}'_p \rho + I_p(1 - \rho)]$, with $0 < \sigma^2 < \infty$ and $-1/(p - 1) < \rho < 1$.

Lemma 1.33 A $p \times p$ compound symmetric covariance matrix may be written

$$\begin{aligned} \Sigma &= \sigma^2 [\mathbf{1}_p \mathbf{1}'_p \rho + I_p(1 - \rho)] \\ &= V \text{Dg}(\lambda) V' \\ &= \begin{bmatrix} \mathbf{v}_0 & \mathbf{V}_t \end{bmatrix} \begin{bmatrix} \lambda_1 & \mathbf{0} \\ \mathbf{0} & \lambda_2 I_{p-1} \end{bmatrix} \begin{bmatrix} \mathbf{v}'_0 \\ \mathbf{V}'_t \end{bmatrix}. \end{aligned} \tag{1.159}$$

The $p \times p$ matrix $\text{Dg}(\lambda) = \text{Dg}(\lambda_1, \lambda_2 \mathbf{1}_{p-1})$ has

$$\lambda_1 = \sigma^2 [1 + (p - 1)\rho] \tag{1.160}$$

$$\lambda_2 = \sigma^2 (1 - \rho). \tag{1.161}$$

The restrictions

$$-1/(p-1) < \rho < 1 \quad (1.162)$$

$$0 < \sigma^2 < \infty \quad (1.163)$$

guarantee $0 < \lambda_j < \infty$ and therefore Σ will be positive definite. The first eigenvector, corresponding to λ_1 , is $\mathbf{v}_0 = \mathbf{1}_p/p^{1/2}$ (normalized to unit length). The first eigenvector may be described as spanning (1) the sum when scaled $\mathbf{1}_p$, (2) the average when scaled $\mathbf{1}_p/p$, or (3) the zero-order polynomial trend when scaled $\mathbf{1}_p/p^{1/2}$. The remaining $p-1$ eigenvectors correspond to λ_2 , which has multiplicity $p-1$. The $p \times (p-1)$ matrix \mathbf{V}_t may be taken to contain first- (linear) through $(p-1)$ -order orthogonal polynomial trend coefficients (scaled to unit length and hence orthonormal). Although \mathbf{v}_0 is unique (up to scaling), \mathbf{V}_t is not. Any other $p \times (p-1)$ orthonormal matrix which is orthogonal to \mathbf{v}_0 will also suffice. In any case, the eigenvectors can always be expressed as a known and constant matrix, no matter what the unknown parameters (ρ, σ^2) .

Proof. Left as an exercise. Hints: verify directly that the eigenvectors reproduce themselves; eigenvalues must be positive, finite and nonzero.

A compound symmetric covariance matrix may be described as having one eigenvalue, λ_1 , of multiplicity 1 with corresponding normalized eigenvector \mathbf{v}_0 , and a second eigenvalue, λ_2 , of multiplicity $p-1$ with corresponding normalized eigenvectors \mathbf{V}_t . Any $p \times (p-1)$ and columnwise orthonormal matrix also orthogonal to \mathbf{u}_0 could be chosen in lieu of the trends.

Definition 1.49 An AR(1), *autoregressive, order 1*, covariance matrix is $\sigma^2 \mathbf{P}$, with $\mathbf{P} = \{\rho_{jk}\}$, with $\rho_{jk} = \rho^{|j-k|}$, for $0 \leq \rho < 1$ and $0 < \sigma^2 < \infty$.

CHAPTER 2

The General Linear Univariate Model

2.1 MOTIVATION

Chapter 2 centers on providing a careful statement of assumptions most often used with the general linear univariate model. A number of specific examples illustrate the basic ideas. Chapters 3, 4, and 5 have the same structure for the multivariate model, multivariate generalizations, and mixed models. Together, applications and properties of the models in Chapters 2–5 will illustrate the need for and uses of the theory in the remainder of the book. Later chapters contain explicit proofs of nearly all basic results in Chapters 2–5. Most others can be deduced with only modest effort from the ones provided.

2.2 MODEL CONCEPTS

Both in the title of and purpose for the book, the univariate, multivariate, and mixed linear models share equal billing. However, the univariate model likely has more pages devoted to it than the multivariate model, which likely has more pages devoted to it than the mixed model. The unevenness reflects the fact that a solid understanding of univariate models allows quickly generalizing many results to multivariate and mixed models. In turn, multivariate theory helps develop mixed model results. The principle holds most often for estimation, while hypothesis testing and inference usually differ in basic ways.

Figure 2.1 illustrates the four aspects of any model: scientific *meaning* of the model, *estimation* of parameters, *inference* about parameters, and *numerical* methods. In Chapters 2–5 we look at the practical interpretations of linear models and focus on scientific meaning. Subsequent groups of chapters center on distribution theory, estimation, inference, and sample size. Although rarely discussed here, accurate numerical methods must be used to ensure the validity of any data analysis.

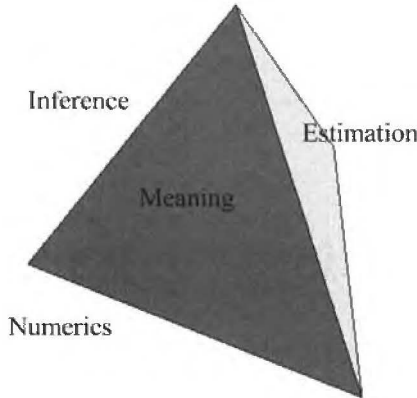


Figure 2.1 Four aspects of a statistical model.

2.3 THE GENERAL LINEAR UNIVARIATE LINEAR MODEL

Definition 2.1 A *general linear univariate model* will be indicated by $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, with *primary parameters* $\{\boldsymbol{\beta}, \sigma^2\}$, and includes the following assumptions.

1. The elements of the $N \times 1$ random vector $\mathbf{y} = \{y_i\}$ are mutually independent.
2. With $\mathbf{X}_i = \text{row}_i(\mathbf{X})$, the $N \times q$ *design matrix*, \mathbf{X} has $\text{rank}(\mathbf{X}) = r \leq q \leq N$, and is fixed and known without appreciable error, conditional on knowing the sampling units, for data analysis. Power analysis requires knowing the predictor distribution in the population.
3. Elements of $\boldsymbol{\beta}$ ($q \times 1$) are fixed and unknown and typically regression coefficients or means.
4. The mean of \mathbf{y} is $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$.
5. Response y_i has finite variance $\sigma^2 \geq 0$, which is fixed and unknown.

Writing $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$ specifies *explicit restrictions* on parameters in $\boldsymbol{\beta}$ through the fixed and known constants \mathbf{R} and \mathbf{a} .

The model is described as *full rank* (FR) if $r = \text{rank}([\mathbf{X}' \mathbf{R}]') = q$ and otherwise as *less than full rank* (LTFR) if $r < q$. Clarity may require writing $\text{GLM}_{N,q}\text{FR}()$ or $\text{GLM}_{N,q}\text{LTFR}()$.

The definition of a GLM describes the “least squares” assumptions because they guarantee estimates for the primary parameters $\boldsymbol{\beta}$ and σ^2 exist which satisfy the least squares criterion. It is very important to recognize that no particular distribution has been specified for any random variable. The rather modest

requirement of finite second (and implicitly first) moments is made. Much more importantly, the requirements of independent and homogeneous observations place strong restrictions on the range of models. The model definition specifies three components: the response for the independent sampling unit, the mean of the response, and the variance of the response.

A number of implications of the GLM definition (least squares assumptions) may be deduced easily. First, response vector \mathbf{y} has finite covariance $\mathcal{V}(\mathbf{y}) = \sigma^2 \mathbf{I}$. Independence makes the off-diagonal elements zero (independence implies zero covariance and correlation), while the homogeneity assumption provides the diagonal elements. Second, the response vector may be separated into purely fixed and purely random vectors by centering the responses to define

$$\mathbf{e} = \mathbf{y} - E(\mathbf{y}) = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}, \tag{2.1}$$

with $E(\mathbf{e}) = \mathbf{0}$ and $\mathcal{V}(\mathbf{e}) = \sigma^2 \mathbf{I}$. Third,

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}. \tag{2.2}$$

The $N \times 1$ constant vector $\mathbf{X}\boldsymbol{\beta}$ describes (models) the first moment of the responses, the mean vector. The $N \times 1$ random vector \mathbf{e} describes (models) the second and higher moments.

Choosing estimators for the primary parameters, $\boldsymbol{\beta}$ and σ^2 , which satisfy a variety of optimal properties which are exact, even in small samples, does not require any particular choice of distribution function for the responses. In contrast, the desire to test hypotheses leads to describing distributions of test statistics. Finding exact distributions for small samples usually requires explicit and particular specification of the distribution of the data.

Data analysts often assume the data have a Gaussian distribution. As detailed in Chapter 8, $\mathbf{y}_1 \sim \mathcal{N}_n(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$ indicates the vector \mathbf{y}_1 ($n \times 1$) has a Gaussian distribution, with finite mean $\boldsymbol{\mu}_1$, finite and positive definite covariance matrix $\boldsymbol{\Sigma}_1$, and density $f_1(\mathbf{y}_1) = (2\pi)^{-n/2} |\boldsymbol{\Sigma}_1|^{-1/2} \exp[-(\mathbf{y}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Sigma}_1^{-1} (\mathbf{y}_1 - \boldsymbol{\mu}_1) / 2]$. In turn, constant \mathbf{T} , $m \times n$ with $m > n$, makes $\mathbf{y}_2 = \mathbf{T}\mathbf{y}_1 \sim \mathcal{SN}_m(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$, singular Gaussian, for $\boldsymbol{\mu}_2 = \mathbf{T}\boldsymbol{\mu}_1$ and $\boldsymbol{\Sigma}_2 = \mathbf{T}\boldsymbol{\Sigma}_1\mathbf{T}'$. The deficient rank of $\boldsymbol{\Sigma}_2$ disallows the existence of a density (the distribution function remains well defined). Writing $(\mathcal{S})\mathcal{N}_m(\boldsymbol{\mu}_2, \boldsymbol{\Sigma}_2)$ indicates the distribution may or may not be singular.

Definition 2.2 Writing $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors indicates $y_i \sim \mathcal{N}_1[\text{row}_i(\mathbf{X})\boldsymbol{\beta}, \sigma^2]$.

Following Kleinbaum, Kupper, Muller, and Nizam (1998) and Muller and Fetterman (2002) the GLM assumptions may be summarized with the mnemonic *HILE Gauss*: homogeneity [$\mathcal{V}(y_i) = \mathcal{V}(y_{i'}) = \sigma^2$], independence ($y_i \perp\!\!\!\perp y_{i'}$ if $i \neq i'$), linearity [$E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$], existence of finite second moments, and, *optionally*, Gaussian observations. The mnemonic groups the least squares assumptions together and separates the distribution assumption.

Adding the assumption of Gaussian errors allows deducing additional properties. With HILE Gauss, essentially all of the assumptions are captured by the statement $\mathbf{y} \sim \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$. Obviously the parameters of the Gaussian distribution coincide with the parameters of the GLM. First and second moments fully characterize both. Also, $\mathbf{e} \sim \mathcal{N}_N(\mathbf{0}, \sigma^2\mathbf{I})$.

The theory of the GLM applies for two apparently disparate classes of applications: regression models and Analysis-of-Variance (ANOVA) models. ANOVA models were developed to test the effects of one or more categorical predictors on a Gaussian response with independent and homogenous errors. Regression models were developed to express a continuous response as a function of one or more continuous predictors. Models with both categorical and continuous predictors fall in between and are best thought of simply as linear models. The underlying theory, for data analysis if not always for power analysis, coincides for all of them.

Example 2.1 Benignus, Muller, Smith, Pieper, and Prah (1990) exposed 74 participants to one of five profiles of carbon monoxide (CO) in the air breathed during the study: Air, Low, Medium, High, or Slow. Accuracy of eye-hand coordination was measured before exposure and four times during exposure. Experience in a series of similar studies, coupled with knowledge of the underlying physical process and measurements, convinced the investigators that the logarithm of an accuracy measure follows a Gaussian distribution. The primary planned analysis compared the five groups on mean change from baseline to time 4.

A GLM for the study may be written in many ways. Muller and Fetterman (2002, Chapter 12) reviewed coding schemes for design matrices. Group sizes were $\{14, 15, 15, 15, 15\}$ for Air, Low, Medium, High, and Slow. A classical LTFR ANOVA coding scheme for the 74 responses, $\{d_i\}$, with $d_i = y_{i4} - y_{i0}$ (y_{i4} is the hour 4 response, y_{i0} is the hour 0 response, the baseline), may be written, with all elements of the design (super-) matrix being conforming vectors,

$$\mathbf{d} = \begin{bmatrix} \mathbf{1}_{14} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_A \\ \alpha_L \\ \alpha_M \\ \alpha_H \\ \alpha_S \end{bmatrix} + \mathbf{e}. \quad (2.3)$$

The parameters are not well defined without adding a constraint. Traditionally, it is assumed, with $\alpha_g \in \{\alpha_A, \alpha_S, \alpha_M, \alpha_H, \alpha_S\}$, that $\sum \alpha_g = 0$. Choosing $\mathbf{R} = [\mathbf{0} \ \mathbf{1}'_5]$ and $\mathbf{a} = 0$ specifies the constraints in the context of a restricted model, with restrictions $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$.

Example 2.2 Deleting the first column of the design matrix in equation (2.3) creates a cell mean coding and a full rank-design matrix. No constraints are needed.

Example 2.3 Deleting any one of columns 2–6 in the design matrix for the first example creates a reference cell coding, also full rank, with the deleted column indicating the reference cell. Chapters 12 and 13 in Muller and Fetterman (2002) describe coding schemes and practical aspects of one-way ANOVA.

Example 2.4 An alternative analysis would use a more general model, the full model in every cell. With $y_{0,j}$ indicating the vector of time 0 (baseline) responses and y_1 the 74×1 vector of responses at time 4, the model is

$$y_1 = \begin{bmatrix} 1_{14} & 0 & 0 & 0 & 0 & y_{0,A} & 0 & 0 & 0 & 0 \\ 0 & 1_{15} & 0 & 0 & 0 & 0 & y_{0,L} & 0 & 0 & 0 \\ 0 & 0 & 1_{15} & 0 & 0 & 0 & 0 & y_{0,M} & 0 & 0 \\ 0 & 0 & 0 & 1_{15} & 0 & 0 & 0 & 0 & y_{0,H} & 0 \\ 0 & 0 & 0 & 0 & 1_{15} & 0 & 0 & 0 & 0 & y_{0,S} \end{bmatrix} \begin{bmatrix} \beta_{0A} \\ \beta_{0L} \\ \beta_{0M} \\ \beta_{0H} \\ \beta_{0S} \\ \beta_{1A} \\ \beta_{1L} \\ \beta_{1M} \\ \beta_{1H} \\ \beta_{1S} \end{bmatrix} + e. \quad (2.4)$$

Chapter 16 in Muller and Fetterman (2002) provides details of practical aspects of coding schemes, estimation, and testing for the full model in every cell. Special cases include ANCOVA and difference scores.

Example 2.5 An Analysis of Covariance (ANCOVA) model may be coded

$$y_1 = \begin{bmatrix} 1_{14} & 0 & 0 & 0 & 0 & y_{0,A} \\ 0 & 1_{15} & 0 & 0 & 0 & y_{0,L} \\ 0 & 0 & 1_{15} & 0 & 0 & y_{0,M} \\ 0 & 0 & 0 & 1_{15} & 0 & y_{0,H} \\ 0 & 0 & 0 & 0 & 1_{15} & y_{0,S} \end{bmatrix} \begin{bmatrix} \beta_{0A} \\ \beta_{0L} \\ \beta_{0M} \\ \beta_{0H} \\ \beta_{0S} \\ \beta_1 \end{bmatrix} + e. \quad (2.5)$$

The model assumes equal slopes and is a special case of the full model in every cell. In turn, analysis of difference scores (as discussed earlier) is a special case of an ANCOVA model because it assumes the common slope is 1.0.

2.4 THE UNIVARIATE GENERAL LINEAR HYPOTHESIS

Definition 2.3 In the $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$, functions of the primary parameters are *secondary parameters*.

Three important examples deserve specific mention because they lie at the heart of the general linear hypothesis (GLH). For \mathbf{C} an $a \times q$ matrix of known constants and $\boldsymbol{\theta}_0$ an $a \times 1$ vector of known constants, $a \times 1$ vectors $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ and $\boldsymbol{\theta} - \boldsymbol{\theta}_0$ are secondary parameters. Furthermore, the covariance matrix of the estimators, namely $\mathcal{V}(\hat{\boldsymbol{\theta}}) = \mathcal{V}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)$, is also a secondary parameter.

The general linear (null) hypothesis may be stated in two equivalent ways:

$$\begin{aligned} H_0 : \mathbf{C}\boldsymbol{\beta} &= \boldsymbol{\theta}_0 \\ H_0 : \boldsymbol{\theta} &= \boldsymbol{\theta}_0. \end{aligned} \quad (2.6)$$

If $a > 1$, then the alternative hypothesis necessarily is $H_A : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$. If $a = 1$ (and θ is a scalar, such as a mean difference), then a one-sided alternative may occasionally be preferred. The notation of a Boolean algebra (as in the following definition) greatly simplifies discussions of inference about hypotheses.

Definition 2.4 (a) A *Boolean algebra* (Weisstein, 2003) is the partial order on subsets (of a collection of sets) which is closed under finite union (OR, \vee), intersection (AND, \wedge) and complementation (NOT, \neg).

(b) Each element defines a *Boolean function*, $\mathbb{B}(\{\})$.

(c) By convention, a two-valued Boolean algebra has values TRUE or FALSE, with $\mathbb{B}(\{\}) = 1$ (TRUE) or $\mathbb{B}(\{\}) = 0$ (FALSE).

Definition 2.5 (a) The *null hypothesis* may be written $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$.

(b) The *alternative hypothesis* may be written $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$.

For linear models, the secondary parameter *noncentrality* characterizes the changes in distributions due to changing hypotheses. As discussed later in the chapter, the test statistic is essentially a noncentrality estimator.

Definition 2.6 (a) For a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with \mathbf{X} fixed and known, the *shift parameter* is

$$\begin{aligned} \delta &= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}']^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0). \end{aligned} \quad (2.7)$$

(b) The *noncentrality parameter* is

$$\omega = \delta / \sigma^2 = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}']^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) / \sigma^2. \quad (2.8)$$

The parameters ω and δ play central roles in the distribution theory of hypothesis tests. The parameter ω is scale free in the sense that multiplying \mathbf{y} by a nonzero constant does not change it. Neither δ nor ω vary under a full-rank transformation of the rows of both \mathbf{C} and $\boldsymbol{\theta}_0$, of the form $\mathbf{C}_T = \mathbf{T}\mathbf{C}$ and $\boldsymbol{\theta}_{0T} = \mathbf{T}\boldsymbol{\theta}_0$, with \mathbf{T} ($a \times a$) of rank a . The null hypothesis is $\boldsymbol{\theta} = \boldsymbol{\theta}_0$, which implies $\delta = (\mathbf{0})' \mathbf{M}^{-1} (\mathbf{0}) = 0$ and therefore $\omega = 0$. The alternative hypothesis has

$\theta \neq \theta_0$, which guarantees $\delta > 0$ and therefore $\omega > 0$. The reader should be cautioned that some authors include a factor of $1/2$ in the definition of ω . The presence or absence of the factor must always be checked in any discussion.

Example 2.6 An independent groups t test with equal sample sizes may be conducted with the following cell mean coding model:

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\ &= \begin{bmatrix} \mathbf{1}_{N/2} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{N/2} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \mathbf{e}. \end{aligned} \tag{2.9}$$

Testing $H_0 : \mu_1 = \mu_2$ leads to using $\mathbf{C} = [1 \quad -1]$ and $\boldsymbol{\theta}_0 = 0$. In turn,

$$\begin{aligned} \mathbf{M} &= \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}' \\ &= [1 \quad -1][[(N/2)\mathbf{I}_2]^{-1}][1 \quad -1]' \\ &= 4/N, \end{aligned} \tag{2.10}$$

$$\begin{aligned} \boldsymbol{\theta} - \boldsymbol{\theta}_0 &= \mathbf{C}\boldsymbol{\beta} - \boldsymbol{\theta}_0 \\ &= [1 \quad -1] \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} - 0 \\ &= \mu_1 - \mu_2, \end{aligned} \tag{2.11}$$

$$\begin{aligned} \delta &= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0) \\ &= (\mu_1 - \mu_2)' (4/N)^{-1} (\mu_1 - \mu_2) \\ &= (\mu_1 - \mu_2)^2 (N/4), \end{aligned} \tag{2.12}$$

and

$$\omega = (N/4)(\mu_1 - \mu_2)^2 / \sigma^2. \tag{2.13}$$

The final form illustrates the principle that noncentrality in the linear model with fixed predictors depends only on sample size, mean differences, and error variance.

A GLH describes a set of linear constraints on $\boldsymbol{\beta}$, namely $\mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$, through the fixed and known constants \mathbf{C} and $\boldsymbol{\theta}_0$. Relative to the unconstrained model, $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, writing $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0, \sigma_0^2)$ specifies a constrained model. Hence every GLH corresponds to comparing a full and constrained model. Any \mathbf{C} matrix with all rows selected from a q -dimensional identity matrix leads to an easily understood constrained model. In such a case, the constrained model may be produced from the full model simply by deleting appropriate columns of \mathbf{X} and corresponding rows of $\boldsymbol{\beta}$. Also, the hypothesis compares the original (full) and the reduced models. The parameter δ has a simple form in the context of comparing the full and constrained models,

$$\delta = N\sigma_0^2 - N\sigma^2. \tag{2.14}$$

2.5 TESTS ABOUT VARIANCES

The general linear hypothesis makes a statement about expected value parameters, elements of β . It is less common, although perfectly reasonable, to consider hypotheses about variances, such as σ^2 . Exact forms for a confidence interval for σ^2 and exact tests of $H_0 : \sigma^2 = \sigma_0^2$ as well as $H_0 : \sigma_1^2 = \sigma_2^2 = \cdots = \sigma_q^2$ are known (with fully independent and Gaussian data).

Tests of variance equality should *not* be used to check the assumption of homogeneity in a linear model. A number of authors have provided analytic and simulation results to support the claim. Instead, robust tests, such as the Satterthwaite approach for the t test, should be used. Similar conclusions apply to testing hypotheses about variances. O'Brien (1979) reviewed available methods and made helpful recommendations.

2.6 THE ROLE OF THE INTERCEPT

The design matrix for the great majority, but not all, linear models either directly includes a column with all 1's, or has columns that span a column of 1's. The following definition formalizes the concept.

Definition 2.7 (a) If constant \mathbf{t}_0 ($q \times 1$) exists such that $\mathbf{X}\mathbf{t}_0 = \mathbf{1}_N$ for $N \times q$ design matrix \mathbf{X} , then the design matrix, and also the associated linear model, *spans an intercept*.

(b) If such a \mathbf{t}_0 exists and $\mathbf{C}\mathbf{t}_0 = \mathbf{0}$, then the hypothesis $H_0 : \mathbf{C}\beta = \theta_0$ *excludes the intercept*.

In the great majority of cases, scientific considerations alone dictate that the model should include an intercept. Most often, but not always, the mean response contains an arbitrary (location) constant and therefore requires the model to span an intercept. Temperature measured in degrees centigrade uses the arbitrary zero point of the temperature at which water freezes. A clinical trial comparing two treatments to reduce fever could compare mean body temperature, with

$$\begin{aligned} H_0 : \mu_1 &= \mu_2 \Leftrightarrow \\ H_0 : (\mu_1 - \mu_2) &= 0. \end{aligned}$$

The parameter $\theta = (\mu_1 - \mu_2)$ does not vary due to adding a constant to the data. Including an intercept allows conducting a location-invariant test. Just as most models span an intercept, most hypotheses exclude the intercept, for the same reasons of scientific irrelevance and lack of meaning.

In some cases, it makes sense to compare the model with and without an intercept. The data for an astronomical study of the temperature of a comet orbiting far from the sun may use values recorded in degrees Kelvin, for which the value zero has meaning. If so, then either excluding an intercept or including one

and testing whether it equals zero may be reasonable. Alternately, Casella (1983) argued that comparing the models with and without an intercept provides a natural way to characterize and understand the statistical value of the intercept.

The preceding brief comments have a number of implications for discussions of the theory of linear models. (1) Most applications will involve models that span an intercept and tests that exclude the intercept. Such tests have location invariance. (2) Completely general results must allow for tests involving the intercept. (3) Completely general results must allow for models (and tests) which do not span an intercept, which corresponds to assuming the intercept equals zero. We urge the reader to always remember the complexities that can occur. Muller and Fetterman (2002, Chapters 4–6) included extensive discussion of the role of the intercept in univariate linear models, tests, and correlation.

2.7 POPULATION CORRELATION AND STRENGTH OF RELATIONSHIP

Data analysts often wish to consider the entire collection of predictors in order to evaluate how well the model fits. Most often, such tests exclude the intercept (for scientific reasons). In the following, β_0 represents the intercept of a GLM, and the intercept may be constrained to be zero, which corresponds to excluding the intercept. In such a model, the special case of a test of all parameters other than the intercept equal to zero compares the original model, $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, to $GLM_{N,1}(y_i; 1 \cdot \beta_0, \sigma_0^2)$ or to $GLM_{N,1}(y_i; 0, \sigma_0^2)$.

Definition 2.8 For the $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, the population value of the *coefficient of determination*, the proportion of variance accounted for by the model equals

$$\rho_u^2 = \frac{N\sigma_0^2 - N\sigma^2}{(N\sigma_0^2 - N\sigma^2) + N\sigma^2} = \frac{\sigma_0^2 - \sigma^2}{\sigma_0^2}. \tag{2.15}$$

For a model spanning an intercept,

$$\rho^2 = \frac{\mathcal{V}(y_i) - \mathcal{V}(y_i|\mathbf{X}_i)}{\mathcal{V}(y_i)} = 1 - \frac{\mathcal{V}(y_i|\mathbf{X}_i)}{\mathcal{V}(y_i)}, \tag{2.16}$$

with ρ the multiple correlation coefficient. In contrast, without an intercept, the general form reduces to

$$\rho_u^2 = \frac{E(y_i^2) - E[(y_i - \mathbf{X}_i\boldsymbol{\beta})^2]}{E(y_i^2)}, \tag{2.17}$$

and ρ_u equals an “uncorrected” (for the intercept) or “generalized” multiple correlation. The parameters ρ^2 and ρ_u^2 may be interpreted as the proportion of

response variance controlled by the predictors (predictable by the model). Both have scale invariance (do not change if the response or any predictors are multiplied by nonzero constants) and ρ^2 also has location invariance (does not change if constants are added to the response or predictors). Kvålseth (1985), Willett and Singer (1988), and Scott and Wild (1991) provided additional useful guidance.

More generally, interest may center on the strength of the relationship between the response variable and some linear combination of predictors, such as a subset, corresponding to a general linear hypothesis, $H_0 : \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$. The resulting parameter equals a squared semipartial multiple correlation. It measures the strength of the relationship between the response and the variables included in the hypothesis, with the predictors in the hypothesis adjusted for predictors not included in the hypothesis. Muller and Fetterman (2002, Chapter 6) discussed correlations in univariate linear models in detail. In the population, the proportion of the response variance controlled by the predictors underlying the general linear hypothesis $H_0 : \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$ may be expressed as

$$\rho_\delta^2 = \frac{\delta}{\delta + N\sigma^2} = \frac{\omega}{\omega + N}. \quad (2.18)$$

Furthermore

$$\omega = \frac{\rho_\delta^2}{(1 - \rho_\delta^2)/N} = N \frac{\rho_\delta^2}{1 - \rho_\delta^2}. \quad (2.19)$$

Consequently, testing the general linear hypothesis $H_0 : \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$ is equivalent to testing $H_0 : \rho_\delta^2 = 0$, or $H_0 : \rho_\delta = 0$, or $H_0 : \omega = 0$, with $\omega = \delta/\sigma^2$. In general, the coefficient of determination and the noncentrality for the hypothesis are one-to-one functions of each other. The coefficient of determination population value does not vary with sample size (while the estimator does). The value falls in the unit interval $0 \leq \rho_\delta^2 \leq 1$, with $\rho_\delta^2 = 0$ [which implies $\mathcal{V}(y_i) = \mathcal{V}(e_i) = \sigma^2 > 0$] occurring only under the null hypothesis and $\rho_\delta^2 = 1$ indicating perfect predictability [and $\mathcal{V}(y_i) > \mathcal{V}(e_i) = \sigma^2 = 0$]. The squared multiple correlation for the entire model arises as a special case by choosing an appropriate \mathbf{C} matrix.

2.8 STATISTICAL ESTIMATES

A less-than-full-rank model has $r = \text{rank}(\mathbf{X}) < q$, and estimation relies on a nonunique generalized inverse $(\mathbf{X}'\mathbf{X})^-$. The unique inverse $(\mathbf{X}'\mathbf{X})^{-1}$ exists only when \mathbf{X} has full rank, $\text{rank}(\mathbf{X}) = q$. With less than full rank,

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{y} \quad (2.20)$$

provides a proper, although biased, estimator of $\boldsymbol{\beta}$, which varies with the choice of generalized inverse. In contrast, a full-rank model defines $\boldsymbol{\beta}$ differently and thereby creates a unique and unbiased estimator, namely

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}. \tag{2.21}$$

Hence a full-rank design defines β in such a way as to guarantee it is *estimable*, in the sense that a known linear function of the data provides an unbiased estimator. The loose definition of “estimable” will be refined in Chapter 11. The obvious estimator of θ may be written

$$\hat{\theta} = \mathbf{C}\tilde{\beta}. \tag{2.22}$$

We delay the discussion as to how to determine whether $\hat{\theta}$ is a good estimator, along with the proofs of all estimation properties, to Chapter 11.

Definition 2.9 For the $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$, (a) *predicted values* are

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{X}\tilde{\beta} \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{H}\mathbf{y}. \end{aligned} \tag{2.23}$$

(b) In turn, the estimated errors, the *residuals*, are

$$\begin{aligned} \hat{\mathbf{e}} &= \mathbf{y} - \mathbf{X}\tilde{\beta} \\ &= [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y} = (\mathbf{I} - \mathbf{H})\mathbf{y}. \end{aligned} \tag{2.24}$$

The projection matrix, $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, earned the nickname “hat matrix” because $\hat{\mathbf{y}} = \mathbf{H}\mathbf{y}$. Very conveniently, \mathbf{H} is unique (Corollary 1.15) even though $(\mathbf{X}'\mathbf{X})^{-1}$ is not, which ensures the predicted values and residuals also are unique.

Definition 2.10 For the $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ and the general linear hypothesis $H_0 : \mathbf{C}\beta = \theta_0$, the *error sum of squares* is $SSE = \hat{\mathbf{e}}'\hat{\mathbf{e}}$ and the *hypothesis sum of squares* is $SSH = \hat{\delta} = (\hat{\theta} - \theta_0)'\mathbf{M}^{-1}(\hat{\theta} - \theta_0)$.

With $r = \text{rank}(\mathbf{X})$ in a univariate model,

$$\begin{aligned} \hat{\sigma}^2 &= \hat{\mathbf{e}}'\hat{\mathbf{e}}/(N - r) \\ &= \mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}/(N - r) \\ &= \mathbf{y}'(\mathbf{I} - \mathbf{H})\mathbf{y}/(N - r) \end{aligned} \tag{2.25}$$

is unique and also can be proven to be an unbiased estimator of σ^2 .

The obvious estimator of noncentrality is

$$\begin{aligned} \hat{\omega} &= \hat{\delta}/\hat{\sigma}^2 \\ &= \frac{(\hat{\theta} - \theta_0)'\mathbf{M}^{-1}(\hat{\theta} - \theta_0)}{\hat{\sigma}^2} \\ &= \frac{SSH}{SSE/(N - r)}. \end{aligned} \tag{2.26}$$

In turn, the usual F statistic is simply $F = \widehat{\omega}/a$, which can be characterized as the average amount of noncentrality per degree of freedom.

The various correlation parameter estimators may be computed in various ways. A general approach may be inferred from

$$\widehat{\rho}_\delta^2 = \frac{SSH}{SSH + SSE}. \quad (2.27)$$

The special case of the usual squared correlation for a model including an intercept (a special case of spanning an intercept) can be computed in the terms of a test of all slopes equal zero. Alternately, for any model spanning an intercept,

$$\widehat{\rho}^2 = 1 - \frac{\mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}}{\mathbf{y}'[\mathbf{I} - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}']\mathbf{y}}. \quad (2.28)$$

The “uncorrected” (for the intercept) or “generalized” multiple correlation estimator assumes \mathbf{X} does not include or span $\mathbf{1}_N$. It takes the form

$$\widehat{\rho}_u^2 = 1 - \frac{\mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}}{\mathbf{y}'\mathbf{y}}. \quad (2.29)$$

2.9 TESTING THE GENERAL LINEAR HYPOTHESIS

Throughout, we assume $F \sim F(\nu_1, \nu_2, \omega)$ indicates the random variable F follows a noncentral F distribution with ν_1 numerator degrees of freedom, ν_2 denominator degrees of freedom, and noncentrality ω , as detailed in Chapter 9. A central case has $\omega = 0$ and is indicated $F \sim F(\nu_1, \nu_2)$.

As detailed in Chapter 15, well-defined tests of *testable* hypotheses require, with $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ of dimensions $a \times a$,

$$\text{rank}(\mathbf{M}) = a, \quad (2.30)$$

which means \mathbf{M} is full rank and invertible, as well as

$$\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X}). \quad (2.31)$$

Only less-than-full-rank models require checking the second condition, which can be replaced by a variety of equivalent conditions (which are usually less computationally convenient). Here we restrict attention to such testable hypotheses. The first condition implicitly requires \mathbf{C} has full (row) rank of a , which provides a sufficient condition to ensure testability for full-rank models but not for less-than-full-rank models.

Definition 2.11 (a) The *test* of hypothesis H_0 is a Boolean function, and also a random variable, $\phi(\mathbf{y}) = \mathbb{B}(\mathbf{y} \in RR)$, and is defined by

(b) *rejection region (critical region)* $RR = \{\mathbf{y} : t > t_0\}$ for

(c) *test statistic* t , a function of the data, \mathbf{y} , and an

(d) appropriate *critical value*, t_0 .

(e) Here $\phi(\mathbf{y}) = 1$ if $t > t_0$ and $\phi(\mathbf{y}) = 0$ if $t \leq t_0$.

(f) The complement of RR is the *acceptance region*, AR .

(g) The *size* of the test is $\alpha = \Pr\{\mathbf{y} \in RR | H_0 = 1\}$ (the null is true).

(h) *Power* is $1 - \beta = \Pr\{\mathbf{y} \in RR | H_A = 1\}$ (the alternative is true).
More generally, the *power function* is $\Pr\{\mathbf{y} \in RR | H\}$.

(i) The probability of making a *type I error* (false positive) is α .

(j) The probability of making a *type II error* (false negative) is β .

The generality of the power function allows discussing both null and alternative hypotheses. If the size of the test (size of RR) is α , then $\Pr\{\mathbf{y} \in RR | H = H_0\} = \alpha$. Furthermore, $\Pr\{\mathbf{y} \in RR | H = H_A\} = 1 - \beta$.

A general linear hypothesis may be tested with the statistic

$$\begin{aligned}
 F &= \frac{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) / a}{\hat{\sigma}^2} \\
 &= \frac{SSH/a}{SSE/(N-r)} = \hat{\omega} / a \\
 &\sim F(a, N-r, \omega).
 \end{aligned}
 \tag{2.32}$$

Under the null, $\omega = 0 \Leftrightarrow \rho_\delta^2 = 0 \Leftrightarrow \delta = 0 \Leftrightarrow \boldsymbol{\theta} = \boldsymbol{\theta}_0$. Except for the simple constant factor of a , the hypothesis degrees of freedom, the usual test statistic merely estimates the noncentrality parameter, ω . The population value (ω) may be interpreted as a times the F statistic which would occur if $\hat{\boldsymbol{\theta}} = \boldsymbol{\theta}$ and $\hat{\sigma}^2 = \sigma^2$.

The test has many optimal properties. It provides the likelihood ratio test and the union-intersection principle test. It always has scale invariance and also has location invariance if the model spans an intercept while the test excludes the intercept. The test has uniformly most power among all similarly invariant tests of size α . Finally, it always has exact size α .

Estimators from $GLM_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$, and $GLM_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta} | \mathbf{C} \boldsymbol{\beta} = \boldsymbol{\theta}_0, \sigma_0^2)$ allow writing an alternate form

$$\hat{\delta} = (N-r+a)\hat{\sigma}_0^2 - (N-r)\hat{\sigma}^2.
 \tag{2.33}$$

In turn,

$$\begin{aligned}
 \hat{\rho}_\delta^2 &= \frac{\hat{\delta}}{\hat{\delta} + (N-r)\hat{\sigma}^2} = \frac{\hat{\omega}}{\hat{\omega} + (N-r)} \\
 &= \frac{aF}{aF + (N-r)} = \frac{SSH}{SSH + SSE},
 \end{aligned}
 \tag{2.34}$$

and

$$\hat{\omega} = \frac{\hat{\rho}_\delta^2}{(1 - \hat{\rho}_\delta^2)/(N - r)} = (N - r) \frac{\hat{\rho}_\delta^2}{1 - \hat{\rho}_\delta^2}. \quad (2.35)$$

Also, the test statistic may be written as

$$F = \frac{\hat{\rho}_\delta^2/a}{(1 - \hat{\rho}_\delta^2)/(N - r)}. \quad (2.36)$$

The forms for parameters and estimators just given prove useful in both developing and understanding discussions of scale and location invariance properties.

2.10 CONFIDENCE REGIONS FOR θ

Confidence intervals, and more general confidence regions, convey useful information about the location of a parameter by combining precision and location properties of parameter estimators. Hypothesis tests also combine information about location and precision. Although tests provide apparently distinct information to that provided by confidence regions, an invertible one-to-one relationship exists between confidence regions and hypothesis tests. Confidence regions can be obtained by inverting hypothesis tests, and a confidence region can be inverted to yield a hypothesis test. In Section 15.6 we explain how the inversions are performed and prove that confidence regions exist only for parameters that are estimable and testable.

Definition 2.12 (a) If data vector \mathbf{y} depends on a primary or secondary parameter vector θ ($a \times 1$) with an unknown true value in parameter space S , then $R(\mathbf{y}) \in S$ with $\Pr\{\theta \in R(\mathbf{y})\} = c(\alpha) \in [0, 1]$ and boundaries defined by a vector-valued function $\mathbf{g}(\mathbf{y}; \alpha)$ is a *confidence region* for θ .

(b) Here $c(\alpha)$ is a *confidence coefficient*.

(c) An exact confidence region has $c(\alpha) = 1 - \alpha$.

(d) An approximate confidence region has $c(\alpha) \approx 1 - \alpha$.

(e) A conservative confidence region has $c(\alpha) \geq 1 - \alpha$.

At least in linear models, $\mathbf{g}(\mathbf{y}, \alpha)$ usually has a closed-form representation. Consequently, computing confidence intervals for parameters of linear models typically proves quite simple.

2.11 SUFFICIENT STATISTICS FOR THE UNIVARIATE MODEL

If the univariate linear model $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian distribution is correct, then the SSCP matrix

$$S = \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{y} \\ \mathbf{y}'\mathbf{X} & \mathbf{y}'\mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{X}' \\ \mathbf{y}' \end{bmatrix} [\mathbf{X} \ \mathbf{y}] \tag{2.37}$$

contains all of the complete sufficient statistics. If \mathbf{X} contains an intercept, then, without loss of generality, \mathbf{X} may be arranged with the intercept in column 1, with $\mathbf{X} = [\mathbf{1}_N \ \mathbf{X}_1]$, for \mathbf{X}_1 of dimension $N \times (q - 1)$. In turn, the SSCP matrix

$$S = \begin{bmatrix} \mathbf{1}'\mathbf{1} & \mathbf{1}'\mathbf{X}_1 & \mathbf{1}'\mathbf{y} \\ \mathbf{X}'_1\mathbf{1} & \mathbf{X}'_1\mathbf{X}_1 & \mathbf{X}'_1\mathbf{y} \\ \mathbf{y}'\mathbf{1} & \mathbf{y}'\mathbf{X}_1 & \mathbf{y}'\mathbf{y} \end{bmatrix} = \begin{bmatrix} \mathbf{1}' \\ \mathbf{X}'_1 \\ \mathbf{y}' \end{bmatrix} [\mathbf{1} \ \mathbf{X}_1 \ \mathbf{y}] \tag{2.38}$$

contains all of the complete sufficient statistics for estimation of all parameters of the regression model identified by the relationship

$$E(\mathbf{y}|\mathbf{X}) . \tag{2.39}$$

Here $[\mathbf{1} \ \mathbf{X}_1 \ \mathbf{y}]$ is $N \times (q + 1)$ so S is $(q + 1) \times (q + 1)$. All parameter estimators and general linear hypothesis tests depend on the data only through the elements of S . Very conveniently, the raw data are not needed for parameter estimation or testing the general linear hypothesis.

EXERCISES

2.1 Provide an explicit interpretation of each parameter in Example 2.5.

2.2 Provide a reference-cell-style coding design matrix for Example 2.5. Include clear specifications of all dimensions. Provide an explicit interpretation of each parameter.

2.3 Provide an effect style coding design matrix for Example 2.6. Include clear specifications of all dimensions. Provide an explicit interpretation of each parameter.

2.4 Explicitly describe the vector of constants \mathbf{t} which demonstrates that the design matrix for Example 2.2 spans an intercept.

2.5 Explicitly describe the vector of constants \mathbf{t} which demonstrates that the design matrix for Example 2.4 spans an intercept.

2.6 For a GLM, briefly explain why multiplying the response values by a nonzero constant automatically implies that $\boldsymbol{\beta}$ and e have also been multiplied by the same nonzero constant for the model to remain valid.

2.7 Prove explicitly that multiplying the model equation by a nonzero constant does not change ω .

2.8 Prove explicitly that multiplying the observed response values by a nonzero constant does not change $\hat{\omega}$ or the observed statistic.

2.9 Give four words or short phrases which highlight the four aspects of a statistical model.

2.10 By default in an ANOVA model, SAS version 8 PROC GLM with the CLASS statement creates a less-than-full-rank design matrix. It creates the design matrix in the following steps: (1) include a column of 1's for the intercept; (2) if a factor has G levels, generate G columns in the design matrix, with each column an indicator variable for one of the levels of the factor. The algorithm (a sweep method; Goodnight, 1979) leads to a “bottom right” reference cell coding scheme, based on the sort order of the formatted values in the class variables (*not* the sort order of the data in the file being analyzed). If factor A has levels $\{1, 2, 3\}$ and factor B has levels $\{x, y, z\}$, the reference cell for a complete two-way design would be 3, z . Additional detail is provided in SAS documentation.

Assume that factor C has four levels $\{1, 2, 3, 4\}$, and it is the only factor in the design (and the class statement). Also let n_i , $i \in \{1, 2, 3, 4\}$, be the number of participants at factor level i and $N = \sum_{i=1}^4 n_i$ be the total number of observations.

2.10.1 Explicitly describe the design matrix originally created by SAS GLM and an associated parameter matrix (give all dimensions and provide brief interpretations). Allow for an unbalanced design (but no missing cells).

2.10.2 Explicitly specify the design matrix and associated parameters implicitly used after the sweep algorithm has been applied.

2.10.3 Explicitly specify constraint matrices and a constrained version of the original model which corresponds to choosing parameters in the reference cell model created by the sweep algorithm.

2.11 Give an example of a nonlinear model which is inherently linear.

2.12 Give an example of a nonlinear model which is not inherently linear.

CHAPTER 3

The General Linear Multivariate Model

3.1 MOTIVATION

The multivariate linear model generalizes the univariate linear model by allowing two or more responses to be measured on each independent sampling unit. Implicitly the model requires that the same design matrix apply to every response and every independent sampling unit have the same set of responses variables. The material in the present chapter summarizes some basic theory needed to fit such models and test hypotheses about predictors and responses.

The multivariate model allows generalizing univariate results for estimation to the multivariate model with very little complication. In contrast, measures of association and hypothesis testing become far more complicated to derive and discuss. The complexity arises from the fact that various criteria for tests lead to a total of nine (yes, nine) commonly used different test methods. Unavoidably, responsible analysis of a multivariate model requires an a priori choice of test method to avoid bias introduced by post hoc p value shopping.

Not surprisingly, the relative appeal of the many tests varies with the nature of the data and the scientific goals of the analysis. Therefore the chapter begins with some suggestions for characterizing dependent responses.

Example 3.1 Example 1.1 contained a classical (LTFR) ANOVA coding for data from Benignus, Muller, Smith, Pieper, and Prah (1990). They exposed 74 human participants to one of five profiles of carbon monoxide (CO) in the air breathed during the study: Air, Low, Medium, High, or Slow. The planned primary analysis compared the five group mean changes from baseline to time 4.

An alternate analysis used four responses, the hour 1–4 differences from baseline. Group sizes were $\{14, 15, 15, 15, 15\}$ for Air, Low, Medium, High, and Slow. The design matrix remains the same as for the univariate model. A multivariate general linear model requires the same design matrix for all responses. LTFR ANOVA coding for the 74 sets of 4 responses, $\{d_{i1}, d_{i2}, d_{i3}, d_{i4}\}$, may be written, with all elements of the design (super-) matrix being conforming vectors,

$$[\mathbf{d}_1 \ \mathbf{d}_2 \ \mathbf{d}_3 \ \mathbf{d}_4] = \begin{bmatrix} \mathbf{1}_{14} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{1}_{15} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \mu_1 & \mu_2 & \mu_3 & \mu_4 \\ \alpha_{A,1} & \alpha_{A,2} & \alpha_{A,3} & \alpha_{A,4} \\ \alpha_{L,1} & \alpha_{L,2} & \alpha_{L,3} & \alpha_{L,4} \\ \alpha_{M,1} & \alpha_{M,2} & \alpha_{M,3} & \alpha_{M,4} \\ \alpha_{H,1} & \alpha_{H,2} & \alpha_{H,3} & \alpha_{H,4} \\ \alpha_{S,1} & \alpha_{S,2} & \alpha_{S,3} & \alpha_{S,4} \end{bmatrix} + \mathbf{E}. \quad (3.1)$$

Here $[\mathbf{d}_1 \ \mathbf{d}_2 \ \mathbf{d}_3 \ \mathbf{d}_4]$ is 74×4 , with columns corresponding to hour and rows to participants. The error matrix has the same dimensions and pattern. Each column in \mathbf{B} contains parameters for a particular hour. The parameters are not well-defined without adding a constraint. For $\alpha_{g,t} \in \{\alpha_A, \alpha_S, \alpha_M, \alpha_H, \alpha_S\}$, it was traditionally assumed $\sum_g \alpha_{g,t} = 0$ holds separately for each value of $t \in \{1, 2, 3, 4\}$. Choosing $\mathbf{R}_x = [0 \ \mathbf{1}'_5]$, $\mathbf{R}_y = \mathbf{I}_4$, and $\mathbf{A} = \mathbf{0}_{1 \times 4}$ implements the constraints in a restricted model with restrictions $\mathbf{R}_x \mathbf{B} \mathbf{R}_y = \mathbf{A}$.

3.2 DEFINITION OF THE MULTIVARIATE MODEL

Definition 3.1 A general linear multivariate model will be indicated by $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$, with primary parameters $\{\mathbf{B}, \Sigma\}$, and includes the following assumptions.

1. The rows of the $N \times p$ random matrix \mathbf{Y} are mutually independent, with $\mathbf{Y}_i = \text{row}_i(\mathbf{Y}) = [y_{i1} \ y_{i2} \ \dots \ y_{ip}]$.
2. With $\mathbf{X}_i = \text{row}_i(\mathbf{X})$, the $N \times q$ design matrix \mathbf{X} has $\text{rank}(\mathbf{X}) = r \leq q \leq N$, and is fixed and known without appreciable error, conditional on knowing the sampling units, for data analysis. Power analysis requires knowing the predictor distribution in the population.
3. Elements of \mathbf{B} ($q \times p$) are fixed and unknown and often regression coefficients or means.
4. The mean of \mathbf{Y} is $E(\mathbf{Y}) = \mathbf{X} \mathbf{B}$.
5. Response matrix \mathbf{Y}_i has finite covariance matrix $\Sigma = \Sigma'$, which is fixed, unknown, and positive definite or positive semidefinite.

Writing $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B} | \mathbf{R}_x \mathbf{B} \mathbf{R}_y = \mathbf{A}, \Sigma)$ specifies explicit restrictions on parameters in \mathbf{B} through the fixed and known constants \mathbf{R}_x , \mathbf{R}_y , and \mathbf{A} .

The model is described as full rank (FR) if $r = \text{rank}([\mathbf{X}' \ \mathbf{R}_x']') = q$ and otherwise as less than full rank (LTFR) if $r < q$. Clarity may require writing $\text{GLM}_{N,p,q} \text{FR}()$ or $\text{GLM}_{N,p,q} \text{LTFR}()$.

The definition of a GLM describes the “least squares” assumptions because they guarantee estimates for the primary parameters \mathbf{B} and Σ exist which satisfy the

least squares criterion. It is very important to recognize that no particular distribution has been specified for any random variable. The rather modest requirement of finite second (and implicitly first) moments is made. Much more importantly, the requirements of independent rows of \mathbf{Y} with common covariance and a common design matrix for every column place strong restrictions on the range of models. The model definition specifies three components: the response matrix for the independent sampling unit, the mean of the response matrix, and the covariance of the response matrix.

A number of implications of the multivariate GLM definition (least squares assumptions) may be deduced easily. The first implication is that elements of response matrix \mathbf{Y} have a finite covariance matrix following a simple pattern. The result may be expressed in terms of either $\text{vec}(\mathbf{Y}')$ or $\text{vec}(\mathbf{Y})$, which are both $Np \times 1$ vectors and are merely permutations of each other. The responses may be decomposed by row or column:

$$\mathbf{Y} = \begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \\ \mathbf{Y}_N \end{bmatrix} = [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \cdots \quad \mathbf{y}_p]. \tag{3.2}$$

Each row of \mathbf{Y} contains the p responses for a single independent sampling unit, while each column contains the N responses for a single variable.

With $p \times N$ \mathbf{Y}' , stacking the N transposed rows (each $p \times 1$) creates an $Np \times 1$ vector,

$$\text{vec}(\mathbf{Y}') = \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_N \end{bmatrix}, \tag{3.3}$$

which has $Np \times Np$ covariance matrix

$$\mathcal{V}[\text{vec}(\mathbf{Y}')] = \mathbf{I}_N \otimes \Sigma = \begin{bmatrix} \Sigma & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \Sigma & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \Sigma \end{bmatrix}. \tag{3.4}$$

The last matrix contains N^2 submatrices, each of them $p \times p$, with a value of either Σ or $\mathbf{0}$. Independence between the N rows of \mathbf{Y} makes the off diagonal matrices $\mathbf{0}$ (independence implies zero covariance and correlation). The homogeneity assumption provides the diagonal matrices.

In contrast, stacking the columns presents exactly the same values in a permuted form. The $Np \times 1$ vector

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_p \end{bmatrix} = \text{vec}(\mathbf{Y}) \quad (3.5)$$

has $Np \times Np$ covariance matrix

$$\begin{aligned} \mathcal{V}[\text{vec}(\mathbf{Y})] &= \boldsymbol{\Sigma} \otimes \mathbf{I}_N \\ &= \begin{bmatrix} \sigma_{11}\mathbf{I} & \sigma_{12}\mathbf{I} & \cdots & \sigma_{1p}\mathbf{I} \\ \sigma_{21}\mathbf{I} & \sigma_{22}\mathbf{I} & \cdots & \sigma_{2p}\mathbf{I} \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{p1}\mathbf{I} & \sigma_{p2}\mathbf{I} & \cdots & \sigma_{pp}\mathbf{I} \end{bmatrix}. \end{aligned} \quad (3.6)$$

In the last equation there are p^2 submatrices, each of them $N \times N$, with the i, j submatrix having the value $\sigma_{ij}\mathbf{I}_N$. We allow heterogeneity and correlation across columns, but not rows of \mathbf{Y} .

A second implication is that the response matrix may be separated into purely fixed and purely random matrices by centering the responses to define

$$\begin{aligned} \mathbf{E} &= \mathbf{Y} - \mathbf{E}(\mathbf{Y}) \\ &= \mathbf{Y} - \mathbf{X}\mathbf{B}. \end{aligned} \quad (3.7)$$

with $\mathbf{E}(\mathbf{E}) = \mathbf{0}$ and $\mathcal{V}[\text{vec}(\mathbf{E}')] = \mathcal{V}[\text{vec}(\mathbf{Y}')] = \mathbf{I}_N \otimes \boldsymbol{\Sigma}$. A closely related third implication of the assumptions is that

$$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}. \quad (3.8)$$

The $N \times p$ constant matrix $\mathbf{X}\mathbf{B}$ describes (models) the first moment of the responses, the mean matrix. The $N \times p$ random matrix \mathbf{E} describes (models) the second and higher moments.

Finding estimators for \mathbf{B} and $\boldsymbol{\Sigma}$ which satisfy a variety of optimal properties in small samples does not require any particular distribution function for the data. In contrast, the desire to test hypotheses leads to describing distributions of test statistics as a function of explicit specification of the data distribution.

Data analysts often assume the errors follow a Gaussian distribution. As detailed in Chapter 8, writing $\mathbf{Y} \sim \mathcal{N}_{n,m}(\mathbf{M}, \boldsymbol{\Xi}, \boldsymbol{\Sigma})$ indicates $\boldsymbol{\Xi}$ and $\boldsymbol{\Sigma}$ are symmetric and positive definite or positive semidefinite, and $\text{vec}(\mathbf{Y}') \sim \mathcal{N}_{n,m}[\text{vec}(\mathbf{M}'), \boldsymbol{\Xi} \otimes \boldsymbol{\Sigma}]$, a direct product matrix Gaussian. A singular distribution may be specified by writing $\mathcal{SN}_{n,m}(\mathbf{M}, \boldsymbol{\Xi}, \boldsymbol{\Sigma})$, while writing $(\mathcal{S})\mathcal{N}_{n,m}(\mathbf{M}, \boldsymbol{\Xi}, \boldsymbol{\Sigma})$ indicates the distribution may or may not be singular. If either $\boldsymbol{\Xi}$ or $\boldsymbol{\Sigma}$ is singular, then so is $\boldsymbol{\Xi} \otimes \boldsymbol{\Sigma}$ because the eigenvalues of the direct product are products of the eigenvalues of the original matrices. Also $\text{rank}(\boldsymbol{\Xi} \otimes \boldsymbol{\Sigma}) = \text{rank}(\boldsymbol{\Xi})\text{rank}(\boldsymbol{\Sigma})$.

Definition 3.2 Writing $GLM_{N,p,q}(Y_i; X_i B | R_x B R_y = A, \Sigma)$ with Gaussian errors indicates $Y_i' \sim \mathcal{N}_p\{\text{row}_i(\mathbf{X})\mathbf{B}', \Sigma\}$.
 Equivalently, $\mathbf{Y} \sim \mathcal{N}_{N,p}(\mathbf{X}\mathbf{B}, \mathbf{I}_N, \Sigma)$.

Following Kleinbaum, Kupper, Muller and Nizam (1998), and Muller and Fetterman (2002), the GLM assumptions may be summarized with the mnemonic *HILE Gauss*: Homogeneity ($\mathcal{V}(Y_i') = \mathcal{V}(Y_{i'}') = \Sigma$), Independence ($Y_i \perp\!\!\!\perp Y_{i'}$ if $i \neq i'$), Linearity ($E(\mathbf{Y}) = \mathbf{X}\mathbf{B}$), existence of finite second moments, and *optionally*, Gaussian observations. The mnemonic groups the least squares assumptions together and separates the distribution assumption.

Adding the assumption of Gaussian errors allows deducing additional properties. With *HILE Gauss*, essentially all of the assumptions are captured by the statement $\mathbf{Y} \sim \mathcal{N}_{N,p}(\mathbf{X}\mathbf{B}, \mathbf{I}_N, \Sigma)$. Obviously the parameters of the Gaussian distribution coincide with the parameters of the GLM. First and second moments fully characterize both. Also, $\mathbf{E} \sim \mathcal{N}_{N,p}(\mathbf{0}, \mathbf{I}_N, \Sigma)$.

As in univariate models, the theory of the multivariate GLM applies for two apparently disparate classes of models: linear regression and ANOVA models. Both are special cases of the multivariate GLM. ANOVA models were developed to test the effects of one or more categorical predictors on two or more Gaussian responses with independent and homogenous errors. Multivariate regression models express a set of continuous responses as a function of one or more continuous predictors. Models with both categorical and continuous predictors fall in between and are best thought of simply as multivariate linear models. The underlying theory, for data analysis, if not always for power analysis, coincides for all of them.

3.3 THE MULTIVARIATE GENERAL LINEAR HYPOTHESIS

Definition 3.3 In the $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$, functions of the primary parameters are *secondary parameters*.

Four important examples deserve specific mention because they lie at the heart of the multivariate general linear hypothesis (GLH). With \mathbf{C} an $a \times q$ matrix of known constants, \mathbf{U} a $p \times b$ matrix of known constants, and Θ_0 an $a \times b$ vector of known constants, $a \times b$ matrices $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ and $\Theta - \Theta_0$ are secondary parameters. Furthermore, the covariance matrix of the estimators, namely $\mathcal{V}[\text{vec}(\hat{\Theta})] = \mathcal{V}[\text{vec}(\hat{\Theta} - \Theta_0)]$, as well as $\Sigma_* = \mathbf{U}'\Sigma\mathbf{U}$, are secondary parameters.

The general linear (null) hypothesis may be stated in two equivalent ways:

$$\begin{aligned} H_0 : CB\mathbf{U} &= \Theta_0 \\ H_0 : \Theta &= \Theta_0. \end{aligned} \quad (3.9)$$

If $\max(a, b) > 1$, then the alternative hypothesis is necessarily $H_A : \Theta \neq \Theta_0$. A one-sided alternative may occasionally be preferred when $\max(a, b) = 1$ (and Θ reduces to θ , a scalar, such as a mean difference).

The C matrix ($a \times q$) defines contrasts *between* groups or levels of predictors by computing linear combinations of coefficients of predictor variables, such as means. The C matrix implicitly computes and explicitly allows testing linear combinations of columns of X , the predictor variables.

The U matrix ($p \times b$) defines contrasts *within* an independent sampling unit (such as person) or level of response character (such as time) by computing linear combinations of coefficients of response variables, such as means. The U matrix implicitly computes and explicitly allows testing linear combinations of columns of Y , the response variables in the model $Y = XB + E$. The parameter Σ provides the $p \times p$ covariance matrix among response variables. The parameter $\Sigma_* = U'\Sigma U$ provides the $b \times b$ covariance matrix among transformed (hypothesis) variables in the model $YU = XB\mathbf{U} + E\mathbf{U}$. With repeated measures in Y , often U contains orthogonal or orthonormal polynomial trend contrasts (linear, quadratic, etc., through order $p - 1$). The zero-order trend corresponds to the mean across the times.

Definition 3.4 For a $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with X fixed and full-rank $M = C(X'X)^{-1}C'$, (a) the *shift parameter* is the $b \times b$ matrix

$$\Delta = (\Theta - \Theta_0)' M^{-1} (\Theta - \Theta_0). \quad (3.10)$$

(b) For full-rank $\Sigma_* = U'\Sigma U$ ($b \times b$), the covariance matrix of transformed (hypothesis) variables, the *noncentrality parameter* is the $b \times b$ matrix

$$\Omega = \Delta \Sigma_*^{-1} = (\Theta - \Theta_0)' M^{-1} (\Theta - \Theta_0) \Sigma_*^{-1}. \quad (3.11)$$

Parameters Δ and Ω play central roles in the distribution theory of multivariate hypothesis tests. The noncentrality parameter Ω is scale free in the sense of corresponding to hypothesis variables standardized to have covariance I_b .

As in the univariate case, both Δ and Ω do not vary under any full-rank transformation of the rows of both C and Θ_0 , of the form $C_T = T_B C$ and $\Theta_{0T} = T_B \Theta_0$, for T_B ($a \times a$) of rank a . On the other hand, both Δ and Ω do vary under simultaneous full-rank transformation of the columns of U and Θ_0 , of the form $U_T = U T_W$ and $\Theta_{0T} = \Theta_0 T_W$, for T_W ($b \times b$) of rank b . However, the eigenvalues of Ω do not vary under such a transformation. As discussed briefly below and in more detail in Chapter 16, the eigenvalues of Ω , in addition to the dimensions of the problem, suffice to fully determine the distribution of the

multivariate test statistics. The multivariate test statistics are functions only of the eigenvalues of $\widehat{\Omega}$ and the dimensions. The statistics do not vary under any full-rank transformation of C or U (and Θ_0) as just discussed.

A wide variety of apparently distinct methods correspond to special cases of the multivariate test statistics, include multivariate ANOVA (MANOVA) and the multivariate approach to repeated measures (MULTIREP). However, tests arising from the univariate approach to repeated measures (UNIREF) differ fundamentally in both origin and many properties. The present chapter contains a brief overview, while Chapter 16 has full details. Both the multivariate (MULTIREP) and UNIREF tests arise from the same estimation theory.

With the spectral decomposition $\Sigma_* = \Upsilon \text{Dg}(\lambda) \Upsilon'$, and $\Delta_\Upsilon = \Upsilon' \Delta \Upsilon$, the distributions of the tests in the UNIREF tests depend on the b eigenvalues of Σ_* , λ , and the diagonal elements (not the eigenvalues) of $\Omega_\Upsilon = \Delta_\Upsilon \text{Dg}(\lambda)^{-1} = \Upsilon' \Delta \Sigma_*^{-1} \Upsilon = \Upsilon' \Omega \Upsilon$. As functions of the eigenvalues of both $\widehat{\Delta}$ and $\widehat{\Sigma}_*$ separately, the corresponding test statistics *do* vary under full rank transformation of the columns of U and Θ_0 , of the form $U_T = UT_W$ and $\Theta_{0T} = \Theta_0 T_W$, with T ($b \times b$) of rank b .

Example 3.2 A multivariate independent groups t test is sometimes described as a two sample Hotelling T^2 test. The setting also corresponds to the special case of a MANOVA test for two groups, which implies $U = I_p$ for p responses. With N_g observations in group $g \in \{A, B\}$, Y and E are $(N_A + N_B) \times p$. A reference cell coding, with $\alpha_j = \mu_{B,j} - \mu_{A,j}$ gives X of dimension $(N_A + N_B) \times 2$ and

$$Y = \begin{bmatrix} \mathbf{1}_{N_A} & \mathbf{0} \\ \mathbf{1}_{N_B} & \mathbf{1}_{N_B} \end{bmatrix} \begin{bmatrix} \mu_{A,1} & \mu_{A,2} & \cdots & \mu_{A,p} \\ \alpha_1 & \alpha_2 & \cdots & \alpha_p \end{bmatrix} + E. \tag{3.12}$$

Therefore, with $C = [0 \ 1]$ and $U = I_p$,

$$\begin{aligned} M &= [0 \ 1] \begin{bmatrix} N_A + N_B & N_B \\ N_B & N_B \end{bmatrix}^{-1} [0 \ 1]' \\ &= [0 \ 1] \begin{bmatrix} 1/N_A & -1/N_A \\ -1/N_A & 1/N_B + 1/N_A \end{bmatrix} [0 \ 1]' \\ &= (N_A + N_B)/(N_A N_B), \end{aligned} \tag{3.13}$$

$$\Sigma_* = I_p' \Sigma I_p = \Sigma \tag{3.14}$$

$$\begin{aligned} \Theta - \Theta_0 &= [0 \ 1] \begin{bmatrix} \mu_{A,1} & \mu_{A,2} & \cdots & \mu_{A,p} \\ \alpha_1 & \alpha_2 & \cdots & \alpha_p \end{bmatrix} I_p - \mathbf{0}_{1 \times p} \\ &= [\alpha_1 \ \alpha_2 \ \cdots \ \alpha_p]. \end{aligned} \tag{3.15}$$

If $A = \{\alpha_j \alpha_j'\}$ then

$$\begin{aligned}\Omega &= [\alpha_1 \alpha_2 \cdots \alpha_p]' [(N_A + N_B)/(N_A N_B)]^{-1} [\alpha_1 \alpha_2 \cdots \alpha_p] \Sigma^{-1} \\ &= (N_A N_B)(N_A + N_B)^{-1} \mathbf{A} \Sigma^{-1}.\end{aligned}\quad (3.16)$$

As in the univariate case, noncentrality depends only on sample size, mean differences, and error variance.

The rank of Ω plays an important role in the distribution theory of multivariate linear models. In general, with $s_* = \text{rank}(\Omega)$, it follows that $0 \leq s_* \leq \min(a, b) = s$, with $s_* = 0$ only when $\Theta = \Theta_0$.

Explaining the origin of the bounds on s_* describes many of the relationships among the dimensions and parameters. The matrices Δ and Σ_*^{-1} are symmetric and $b \times b$. A well-defined test requires full (row) rank of a for \mathbf{C} . It also requires full (column) rank of b for \mathbf{U} , which ensures full rank of b for $\Sigma_* = \mathbf{U}' \Sigma \mathbf{U}$ (given a full-rank Σ , as nearly always assumed) and also Σ_*^{-1} . The $b \times b$ noncentrality matrix, $\Omega = \Delta \Sigma_*^{-1}$, will not be symmetric, except for special cases. Furthermore, $\text{rank}(\Omega) = \text{rank}(\Delta)$, with $\Delta = (\Theta - \Theta_0)' \mathbf{M}^{-1} (\Theta - \Theta_0)$. The $a \times a$ \mathbf{M} matrix must be full rank for a well-defined test, which ensures the rank of Δ ($b \times b$) equals the rank of $(\Theta - \Theta_0)$ ($a \times b$). Obviously $0 \leq \text{rank}(\Theta - \Theta_0) \leq \min(a, b)$. Here $\text{rank}(\Theta - \Theta_0) = 0$ only when $\Theta = \Theta_0$.

Example 3.3 The matrix $\Omega = [(N_A N_B)/(N_A + N_B)] \mathbf{A} \Sigma^{-1}$ from the last example has rank 1, which can be proven as follows. Square and full rank Σ^{-1} ensures $\text{rank}(\Omega) = \text{rank}(\mathbf{A})$. If $\alpha = [\alpha_1 \alpha_2 \cdots \alpha_p]'$, then $\mathbf{A} = \alpha \alpha'$. Hence $\text{rank}(\mathbf{A}) = \text{rank}(\alpha \alpha') = \text{rank}(\alpha) = 1$.

Alternately, if $\alpha_1 = \alpha / \sqrt{\alpha' \alpha}$, then there exists orthonormal \mathbf{A}_0 of dimension $p \times (p-1)$ with $\mathbf{A}_0' \alpha_1 = \mathbf{0}$. The spectral decomposition of \mathbf{A} is

$$\begin{aligned}\mathbf{A} &= \alpha \alpha' = \frac{1}{\sqrt{\alpha' \alpha}} \alpha (\alpha' \alpha) \alpha' \frac{1}{\sqrt{\alpha' \alpha}} \\ &= \alpha_1 (\alpha' \alpha) \alpha_1' \\ &= [\alpha_1 \ \mathbf{A}_0] \begin{bmatrix} (\alpha' \alpha) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \alpha_1' \\ \mathbf{A}_0' \end{bmatrix}.\end{aligned}\quad (3.17)$$

The eigenvalues of $\Omega = (N_A N_B)(N_A + N_B)^{-1} \alpha \alpha' \Sigma^{-1}$ and $\Omega_1 = (N_A N_B)(N_A + N_B)^{-1} \alpha' \Sigma^{-1} \alpha$ coincide, but not the eigenvectors. Writing $\omega_1 = (N_A N_B)(N_A + N_B)^{-1} \alpha' \Sigma^{-1} \alpha$ is more precise because ω_1 is 1×1 , a scalar. Any scalar has a single eigenvalue, the scalar itself. The expression may be written

$$\omega_1 = \frac{N_A N_B}{N_A + N_B} [\mu_{B,1} - \mu_{A,1} \cdots \mu_{B,p} - \mu_{A,p}] \Sigma^{-1} \begin{bmatrix} \mu_{B,1} - \mu_{A,1} \\ \vdots \\ \mu_{B,p} - \mu_{A,p} \end{bmatrix}.\quad (3.18)$$

The parameter ω_1 is the Mahalanobis distance from the origin of the vector of group differences. If the responses are uncorrelated, then Σ is diagonal and

$$\begin{aligned} \omega_1 &= \frac{N_A N_B}{N_A + N_B} \alpha' [\text{Dg}(\{\sigma_j^2\})]^{-1} \alpha \\ &= \frac{N_A N_B}{N_A + N_B} \sum_{j=1}^p \frac{(\mu_{B,j} - \mu_{A,j})^2}{\sigma_j^2}. \end{aligned} \tag{3.19}$$

If $p = 1$ and $N_A = N_B$, then $(N_A N_B)(N_A + N_B)^{-1} = N/4$ and ω_1 reduces to the t -test result in Chapter 2.

Writing $\text{GLM}_{N,p,q}(Y_i; X_i B, \Sigma)$ specifies an unconstrained model. A GLH (general linear hypothesis) describes a set of linear constraints on B , namely $C B U = \Theta_0$, through the fixed and known constants C , U , and Θ_0 . Writing $\text{GLM}_{N,p,q}(Y_i; X_i B | C B U = \Theta_0, \Sigma)$ specifies a constrained model. Hence every GLH corresponds to comparing a full and constrained model. Choosing $U = I_p$ has the same effect as not using any U matrix. With $U = I_p$, any C with all rows selected from I_q leads to an easily understood constrained model. In such a case, the constrained model may be produced from the full model simply by deleting appropriate columns of X and corresponding rows of B . Similarly, with $C = I_q$, any U matrix with all columns selected from I_p leads to an easily understood constrained model. In such a case, the constrained model may be produced from the full model simply by deleting appropriate columns of Y and corresponding columns of B . More generally, U defines a priori linear transformations of the response variables, such as an average or set of difference scores. The hypothesis compares the original (full) and the reduced models.

3.4 TESTS ABOUT COVARIANCE MATRICES

The general linear hypothesis makes a statement about expected value parameters, elements of β . It is less common, although perfectly reasonable, to consider hypotheses about covariance matrices, such as Σ or Σ_* . Some exact tests have been developed. Morrison (1990), Anderson (2004), and Timm (2002, Section 3.8) have additional details.

Tests of covariance pattern should *not* be used to check the assumption of homogeneity in a linear model. A number of authors have provided analytic and simulation results to support the proposition. Chapter 16 includes discussion of tests about covariance matrices. The tests do perform appropriately when used as they were intended.

3.5 POPULATION CORRELATION

The univariate correlation applies to one response and one predictor. Recalling $\hat{y}_i = \beta_0 + \beta_1 x_i$ implies that the squared correlation of the response and predictor is the same as that between the response and the predicted value, $\rho^2(y_i, x_i) = \rho^2(y_i, \hat{y}_i)$. Allowing many predictors gives the squared multiple

correlation $\rho^2(y_i, \{x_{i1}, \dots, x_{iq-1}\}) = \rho^2(y_i, \hat{y}_i)$ with $\hat{y}_i = \hat{\beta}_0 + \sum_{j=1}^{q-1} \hat{\beta}_j x_{ij}$ for a model with an intercept [more generally, $\hat{y}_i = \mathbf{X}_i(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$]. Generalizing to p responses and q predictors (which span an intercept) leads to considering a set of $s = \min(p, q - 1)$ squared canonical correlations. Muller (1982) described the model underlying the traditional approach to canonical correlation as a measure of association between two sets of variables. Canonical correlation also lie at the heart of the multivariate linear model.

Definition 3.5 For the $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ and associated general linear hypothesis $H_0 : \mathbf{C}\mathbf{B}\mathbf{U} = \Theta_0$, the (generalized) squared canonical correlations $\{\rho_{*k}^2\}$ are the eigenvalues of the $b \times b$ matrix $(\Delta + N\Sigma_*)(N\Sigma_*)^{-1}$.

The eigenvalues of $\Omega = \Delta\Sigma_*^{-1}$ are $\{\omega_k\}$, with $\omega_k = N\rho_{*k}^2/(1 - \rho_{*k}^2)$. With $\Sigma_* = \Phi\Phi'$, the eigenvalues of Ω coincide with the eigenvalues of the symmetric and positive definite or positive semidefinite and $b \times b$ matrix $\Omega_\Phi = \Phi^{-1}\Delta\Phi^{-t}$. In all cases s_* nonzero values of ρ_{*k}^2 occur, with $s_* = \text{rank}(\Omega_\Phi) = \text{rank}(\Omega) = \text{rank}(\Delta) = \text{rank}(\Theta - \Theta_0)$. Also, $0 \leq s_* \leq s = \min(a, b)$. Only under the null hypothesis does $s_* = 0$. In general,

$$\omega_k = N \frac{\rho_{*k}^2}{1 - \rho_{*k}^2} \tag{3.20}$$

and

$$\rho_{*k}^2 = \frac{\omega_k}{\omega_k + N}. \tag{3.21}$$

Naturally the multivariate formulas reduce to the univariate formulas if $p = 1$. In turn, $b = 1$ implies $s = 1$, $s_* = 1$ under the alternative, and $s_* = 0$ under the null.

As in the univariate case, $\{\omega_k\}$ and $\{\rho_{*k}^2\}$ always have scale invariance. If the model spans an intercept and the hypothesis excludes it, the hypothesis also has location invariance. If so, the generalized canonical correlations are appropriately described simply as canonical correlations.

3.6 STATISTICAL ESTIMATES

In practice, generalizing formulas for estimates from univariate to multivariate models simply requires replacing the $N \times 1$ vectors \mathbf{y} and \mathbf{e} by the $N \times p$ matrices \mathbf{Y} and \mathbf{E} , with the additional columns corresponding to the p responses. Similarly, $N \times 1 \hat{\mathbf{y}}$ and $\hat{\mathbf{e}}$ become $N \times p \hat{\mathbf{Y}}$ and $\hat{\mathbf{E}}$, while $q \times 1 \boldsymbol{\beta}$ and $\hat{\boldsymbol{\beta}}$ become $q \times p \mathbf{B}$ and $\hat{\mathbf{B}}$. The reason lies hidden in three implicit assumptions: (1) the same design matrix applies to all response variables, which correspond to columns of \mathbf{Y} ; (2) there are no missing data; and (3) each value within a variable (column of \mathbf{Y}) was measured in a consistent way (no appreciable mistiming is allowed).

Adding responses adds columns in \mathbf{Y} but does not change the properties of the design matrix. In a LTFR model [$r = \text{rank}(\mathbf{X}) < q$] estimation relies on a nonunique generalized inverse $(\mathbf{X}'\mathbf{X})^-$. In contrast, $(\mathbf{X}'\mathbf{X})^{-1}$ exists only when \mathbf{X} has full rank [$\text{rank}(\mathbf{X}) = q$]. With less than full rank, the matrix

$$\tilde{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{Y} \tag{3.22}$$

is a proper estimator of \mathbf{B} . However, it is biased and varies with the choice of generalized inverse. In contrast, a full-rank model defines \mathbf{B} differently and thereby makes possible a unique and unbiased estimator, namely

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}. \tag{3.23}$$

A full rank design defines \mathbf{B} so that it is *estimable*, which means that a known function of the data is an unbiased estimator. The definition will be refined later.

Definition 3.6 For the $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$,

(a) *predicted values* may be expressed as

$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{X}\tilde{\mathbf{B}} \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{Y} \\ &= \mathbf{H}\mathbf{Y}. \end{aligned} \tag{3.24}$$

(b) In turn, the estimated errors, the *residuals*, are

$$\begin{aligned} \hat{\mathbf{E}} &= \mathbf{Y} - \mathbf{X}\tilde{\mathbf{B}} \\ &= [\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}']\mathbf{Y} \\ &= (\mathbf{I} - \mathbf{H})\mathbf{Y}. \end{aligned} \tag{3.25}$$

The “hat matrix” \mathbf{H} is unique even though $(\mathbf{X}'\mathbf{X})^-$ is not, which also makes the predicted values and residuals unique.

With \mathbf{C} an $a \times q$ matrix of known constants and \mathbf{U} a $p \times b$ matrix of known constants, $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ defines an $a \times b$ matrix of secondary parameters, with corresponding estimator $\hat{\Theta} = \mathbf{C}\hat{\mathbf{B}}\mathbf{U}$ or $\mathbf{C}\hat{\mathbf{B}}\mathbf{U}$. The covariance matrix of the estimator may also be described as a secondary parameter. We delay the discussion as to how to determine whether $\hat{\Theta}$ is a good estimator, along with the proofs of all estimation properties, to Chapter 12.

Definition 3.7 For the $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ and associated general linear

hypothesis $H_0 : \mathbf{C}\mathbf{B}\mathbf{U} = \Theta_0$, the *error sum of squares* is $S_e = \mathbf{U}'\hat{\mathbf{E}}'\hat{\mathbf{E}}\mathbf{U}$, and the *hypothesis sum of squares* is $S_h = \hat{\Delta} = (\hat{\Theta} - \Theta_0)'\mathbf{M}^{-1}(\hat{\Theta} - \Theta_0)$, which are both $b \times b$ matrices. Choosing $\mathbf{U} = \mathbf{I}_p$ gives $S_e = \hat{\mathbf{E}}'\hat{\mathbf{E}}$, the error sums squares for the model.

The sum of squared errors $SSE = \hat{e}'\hat{e}$ generalizes from a scalar to a $b \times b$ matrix $U'\hat{E}'\hat{E}U$, and scalar $\hat{\sigma}^2$ becomes the $b \times b$ matrix $\hat{\Sigma}_* = U'\hat{\Sigma}U$. Row and column labels of $U'\hat{E}'\hat{E}U$ and $\hat{\Sigma}_*$ are the names of transformed response variables (columns of YU , EU , $\hat{Y}U$, and $\hat{E}U$). As with $p = 1$ (the univariate model), $r = \text{rank}(\mathbf{X})$ and

$$\begin{aligned}\hat{\Sigma} &= \hat{E}'\hat{E}/(N-r) \\ &= Y'[I - X(X'X)^{-1}X']Y/(N-r) \\ &= Y'(I - H)Y/(N-r)\end{aligned}\tag{3.26}$$

is unique and also can be proven to be an unbiased estimator of Σ . Furthermore

$$\hat{\Sigma}_* = U'Y'(I - H)YU/(N-r)\tag{3.27}$$

is also unique and unbiased. With $p = 1$, the univariate model, necessarily $U = [1] = I_1$. Multivariate models allow for $U \neq I_p$, which leads to interest in $\Sigma_* = U'\Sigma U$.

The obvious estimators are $\hat{\Delta} = (\hat{\Theta} - \Theta_0)'M^{-1}(\hat{\Theta} - \Theta_0) = S_h$ and

$$\hat{\Omega} = \hat{\Delta}\hat{\Sigma}_*^{-1} = (\hat{\Theta} - \Theta_0)'M^{-1}(\hat{\Theta} - \Theta_0)\hat{\Sigma}_*^{-1}.\tag{3.28}$$

3.7 OVERVIEW OF TESTING MULTIVARIATE HYPOTHESES

As detailed in Chapter 16, well-defined multivariate tests of *testable* hypotheses require three properties. With $M = C(X'X)^{-1}C'$ of dimensions $a \times a$,

$$\text{rank}(M) = a,\tag{3.29}$$

which means M is full rank and invertible. Testable hypotheses require estimable parameters, which is guaranteed when

$$C = C(X'X)^{-1}(X'X).\tag{3.30}$$

Ensuring a testable multivariate hypothesis imposes an additional requirement, namely full (column) rank of the $p \times b$ matrix U :

$$\text{rank}(U) = b,\tag{3.31}$$

A testable hypothesis necessarily has $b \leq p$. Only less-than-full-rank models require checking the estimability condition, which can be replaced by a variety of equivalent conditions (which are usually less computationally convenient). The first condition implicitly requires C , of dimension $a \times q$, have full (row) rank a , which provides a sufficient condition to ensure testability for full-rank models, but not for less-than-full-rank models. In any case, a testable hypothesis necessarily has $a \leq q$.

A multivariate set of response variables may be tested with any one of nine different tests, which fall into three groups, as summarized in Table 3.1. In

contrast to the univariate setting, no single test satisfies the various measures of goodness for the multivariate general linear hypothesis $H_0 : CBU = \Theta_0$ for all conditions. Although the Wilks statistic (WLK) provides the likelihood ratio test, Roy's largest root (RLR) provides the union-intersection principle test. Equally important, the relative powers of the various tests vary with the pattern of noncentralities, which implies each may be preferred in some settings. We shall provide more detail after providing explicit forms for the tests.

Table 3.1 Tests for Multivariate Hypotheses

Approach	Test	Σ_* Eigenvalues	Test Size	Best Power?
Bonferroni		Any	$\leq \alpha$	$\Sigma = \text{Dg}(\{\sigma_j^2\})?$
MULTIREP (MANOVA)	Hotelling-Lawley (HLT)	Any	$= \alpha$	$s_* > 1$
	Pillai-Bartlett (PBT)	Any	$= \alpha$	$s_* > 1$
	Wilks Likelihood (WLK)	Any	$= \alpha$	$s_* > 1$
	Roy's Largest Root (RLR)	Any	$= \alpha$	$s_* = 1$
UNIREP	Box Conservative	Any	$\leq \alpha$	$\epsilon = 1/b$
	Geisser-Greenhouse (GG)	Any	$\lesssim \alpha$	ϵ near 1
	Huynh-Feldt (HF)	Any	$\approx \alpha$	ϵ near 1
	Uncorrected (UN)	Any	$\geq \alpha$	does not apply
		$\epsilon = 1$	$= \alpha$	UMP ¹

¹Uniformly most powerful given assumptions

With p variables, the simplest approach uses a Bonferroni correction with α/p test size in separate univariate analyses of each response variable. The approach does not completely test the multivariate general linear hypothesis $H_0 : CBU = \Theta_0$. Instead, for each C (between subject, a group contrast), a set of p distinct U matrices are created, with each a distinct column of I_p . The approach has two particular strengths. First, it tolerates missing data and different design matrices for each response. Second, it allows allocating more test size to the important variables, which may better meet the scientific goals and desires of the investigators. The approach obviously inherits the invariance properties of the univariate model.

The second approach, labeled MULTIREP, groups four tests together, traditionally described as the multivariate tests because they were specifically developed for the multivariate general linear hypothesis $H_0 : CBU = \Theta_0$. They share the important property of invariance to the value of Σ and therefore to Σ_* . The tests may be employed with the "multivariate" approach to repeated measures, in contrast to the "univariate" approach to repeated measures (UNIREP). Both approaches to repeated measures decompose the p dimensional response space into

the $p - 1$ dimensions within- plus 1 dimension between-“subject” (independent sampling unit) subspaces. Tests involving contrasts between subject use $\mathbf{u}_B = \mathbf{1}_p/p$, of dimension $p \times 1$, which computes the average response across repeated measures. Tests involving contrasts within subject's use \mathbf{U}_W of dimension $p \times (p - 1)$, with full column rank of $p - 1$, such that $\mathbf{I}_p = [\mathbf{u}_B \ \mathbf{U}_W] \mathbf{T}$, with \mathbf{T} full rank and $p \times p$. Two common choices include $\mathbf{U}_W = [\mathbf{1}_{p-1} \ -\mathbf{I}_{p-1}]'$, which gives pairwise contrasts, and the $p \times (p - 1)$ matrix of orthonormal polynomial trend coefficients, with columns corresponding to linear, quadratic, cubic, etc. Values of the continuous variable labeling the columns of \mathbf{Y} , such as time, provide the points needed to generate the coefficients.

The four MULTIREP tests also provide an appropriate test of a MANOVA hypothesis, which always uses $\mathbf{U} = \mathbf{I}_p$. The fact that $\mathbf{I}_p = [\mathbf{u}_B \ \mathbf{U}_W] \mathbf{T}$, with \mathbf{T} full rank and $p \times p$ allows concluding that the MANOVA hypothesis spans both the between and within hypothesis contrasts. It asks whether *any* linear combination of the responses included in the hypothesis (by the choice of \mathbf{C}) is related to the predictors.

The uncorrected UNIREP test was developed long before any of the other UNIREP tests. Validity of the uncorrected test requires $\Sigma_* = \mathbf{U}' \Sigma \mathbf{U} = \mathbf{I}_b \sigma_*^2$, namely sphericity of the b transformed responses corresponding to the general linear hypothesis. With sphericity the uncorrected test provides an exact size- α test with uniformly most power among all similarly invariant tests. Without sphericity the uncorrected test can have greatly inflated test size.

The covariance matrix of the original responses achieves compound symmetry when all p response have the same variance, σ^2 , and all $p(p - 1)/2$ pairs of distinct responses have the same correlation, ρ (Lemma 1.33 summarizes properties). Choosing \mathbf{U} to be either $\mathbf{u}_0 = \mathbf{1}_p/p^{1/2}$ or any \mathbf{U}_t with (1) $\mathbf{U}_t' \mathbf{1}_p = \mathbf{0}$ and (2) $\mathbf{U}_t' \mathbf{U}_t = \mathbf{I}_b$ combines with compound symmetry to provide a sufficient (but not necessary) set of conditions to guarantee sphericity. If $\mathbf{U} = \mathbf{u}_0$, then $b = 1$ and $\mathbf{U}' \Sigma_{CS} \mathbf{U} = \lambda_1 = \sigma^2[1 + (p - 1)\rho]$. With repeated measures, the data may be arranged, without loss of generality, such that all observations in column j were collected at time t_j , with $t_j \in \{t_1, t_2, \dots, t_p\}$ and $t_j < t_{j+1}$. Here and throughout the book, time may be thought of as a metamer representing any interval- or ratio-scale dimension along which the observations vary within a subject (independent sampling unit, ISU). The orthonormal polynomial trends provide one convenient choice for \mathbf{U}_t .

Example 3.4 If $p = 3$ and the times are equally spaced, the orthonormal trends may be taken to be

$$\mathbf{U}_t = \begin{bmatrix} -1 & 1 \\ 0 & -2 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 0 \\ 0 & 1/\sqrt{6} \end{bmatrix}. \tag{3.32}$$

The first column provides the linear trend and the second column provides the quadratic trend. In turn, $\mathbf{U}_t' \Sigma_{CS} \mathbf{U}_t = \lambda_2 \mathbf{I}_{p-1} = [\sigma^2(1 - \rho)] \mathbf{I}_{p-1}$. Combining the

forms gives the spectral decomposition

$$\Sigma_{CS} = [\mathbf{u}_0 \mathbf{U}_t] \begin{bmatrix} \lambda_1 & \mathbf{0} \\ \mathbf{0} & \lambda_2 \mathbf{I}_{p-1} \end{bmatrix} \begin{bmatrix} \mathbf{u}'_0 \\ \mathbf{U}'_t \end{bmatrix}. \tag{3.33}$$

Without sphericity and with λ_k indicating one of the b possibly distinct eigenvalues of Σ_* , Box (1954a, b) proposed quantifying the deviation from sphericity with the parameter

$$\begin{aligned} \epsilon &= \frac{\text{tr}^2(\Sigma_*)}{b \text{tr}(\Sigma_*^2)} \\ &= \frac{\left(\sum_{k=1}^b \lambda_j/b\right)^2}{\left(\sum_{k=1}^b \lambda_j^2\right)/b} = \frac{(\bar{\lambda})^2}{\lambda^2}. \end{aligned} \tag{3.34}$$

Under sphericity $\epsilon = 1$, while in general $1/b \leq \epsilon \leq 1$. Test size inflation for $\epsilon < 1$ led to the development of the Box conservative, Geisser-Greenhouse, and Huynh-Feldt tests. All use the same test statistic as the uncorrected test but use distinct and more stringent critical values.

Allowing for any covariance pattern and for $s = \min(a, b) > 1$, no single test provides the uniformly most powerful test (among size α and similarly invariant tests). The special case $b = 1$ implies $s = \min(a, b) = 1$ and causes all MULTIREP and UNIREP tests to become equivalent by providing exactly the same F statistic and p value. Furthermore, if $b = 1$ the single test provides the usual exact size- α and uniformly most powerful test (among the class of unbiased and scale invariant tests). Furthermore, $b > 1$ and $a = 1$ imply $s = \min(a, b) = 1$. If $a = 1$, then the MULTIREP tests become equivalent to each other by providing exactly the same F statistic and p value. Similarly, if $a = 1$, then the (single) MULTIREP test provides an exact size- α and uniformly most powerful test among the class of unbiased and scale invariant tests. However, if $a = 1$ and $b > 1$ the UNIREP tests are not equivalent to each other or to the MULTIREP test. The differences arise from the fact that the UNIREP tests are only invariant to an orthonormal transformation, rather than fully scale invariant. Both UNIREP and MULTIREP tests correspond to transforming the model and hypothesis as follows:

$$\begin{aligned} \mathbf{YU} &= \mathbf{XBU} + \mathbf{EU} \\ \mathbf{Y}_u &= \mathbf{XB}_u + \mathbf{E}_u, \end{aligned} \tag{3.35}$$

\mathbf{E}_u has $b \times b$ covariance $\Sigma_* = \mathbf{U}'\Sigma\mathbf{U}$,

$$\begin{aligned} H_0 &: \mathbf{CBU} = \Theta_0 \\ H_0 &: \mathbf{CB}_u \mathbf{I}_b = \Theta_0. \end{aligned} \tag{3.36}$$

With spectral decomposition $\Sigma_* = \mathbf{YDg}(\lambda)\mathbf{Y}'$ and $\mathbf{Y}'\mathbf{Y} = \mathbf{I}_b$ the MULTIREP and UNIREP tests do not vary under orthonormal transformation of the model and Θ_0 . In particular

$$\begin{aligned} Y_u \Upsilon &= X B_u \Upsilon + E_u \Upsilon \\ Y_{\Upsilon} &= X B_{\Upsilon} + E_{\Upsilon} \end{aligned} \tag{3.37}$$

and $\Theta_{0\Upsilon} = \Theta_0 \Upsilon$. Here E_{Υ} has $b \times b$ covariance $\Sigma_{\Upsilon} = \text{Dg}(\lambda)$. If $\Phi = \Upsilon \text{Dg}(\lambda)^{1/2}$, then $\Phi^{-t} = \Upsilon \text{Dg}(\lambda)^{-1/2}$. The eigenvector transformation eliminates all correlations. With Gaussian data, the process creates independent but typically heterogeneous responses. In contrast to the situation for orthonormal invariance, in general the MULTIREP tests do *not* vary while the UNIREP tests *do* vary under the (particular scale) transformation

$$\begin{aligned} Y_{\Upsilon} \text{Dg}(\lambda)^{-1/2} &= X B_{\Upsilon} \text{Dg}(\lambda)^{-1/2} + E_{\Upsilon} \text{Dg}(\lambda)^{-1/2} \\ Y_u \Phi^{-t} &= X B_u \Phi^{-t} + E_u \Phi^{-t} \\ Y_{\Phi} &= X B_{\Phi} + E_{\Phi}. \end{aligned} \tag{3.38}$$

Here E_{Φ} has $b \times b$ covariance $\Sigma_{\Phi} = I_b$. The only exception occurs when $\Sigma_* = \sigma_*^2 I_b$, namely sphericity, with $\epsilon = 1$, and then the uncorrected test achieves size α and uniformly more power than the MULTIREP test. If $a = 1, b > 1$, and $\epsilon < 1$, examples can be found in which either the MULTIREP test (exact size α) or a corrected UNIREP test (at least approximately size α) can be more powerful than the other for a particular Σ_* . The same statement holds for $a > 1, b > 1$ [$s = \min(a, b) > 1$], and $\epsilon < 1$.

3.8 COMPUTING MULTIREP TESTS

The four MULTIREP test statistics are all simple functions of the $s = \min(a, b)$ nonzero estimators of the generalized squared canonical correlations $\{\hat{\rho}_{*k}^2\}$. The correlations are generalized in the sense that they may or may not be adjusted for an intercept. The $\{\hat{\rho}_{*k}^2\}$ are the nonzero eigenvalues of the $b \times b$ matrix $\widehat{\Delta}[\widehat{\Delta} + (N - r)\widehat{\Sigma}_*]^{-1}$, with $r = \text{rank}(X)$.

Definition 3.8 A measure of *multivariate association* generalizes the concept of a squared multiple correlation, the proportion of variance controlled by the hypothesis. Each of the multivariate test statistics leads to a different measure of association. For test statistic m , $0 \leq \eta_m \leq 1$, with $\eta_m = 0$ corresponding to no relationship and $\eta_m = 1$ to a perfect relationship.

Table 3.2 gives expressions for each of the MULTIREP test statistics. For compactness and to emphasize the sums-of-squares nature of the matrices, the alternate notations $\nu_e = N - \text{rank}(X)$, $S_h = \widehat{\Delta}$, and $S_e = (N - r)\widehat{\Sigma}_*$ are used. The column labeled $\hat{\eta}_m$ defines a measure of the strength of multivariate association for test m which corresponds to the form of the test statistic. Cramer and Nicewander (1979) reviewed measures of multivariate association in the context of two models, $Y = X B_{Y|X} + E_{Y|X}$ and $X = Y B_{X|Y} + E_{X|Y}$. Such

measures of correlation all share the property of symmetry with respect to the two sets of variables \mathbf{X} and \mathbf{Y} in the GLM: $\eta_m(\mathbf{Y}, \mathbf{X}) = \eta_m(\mathbf{X}, \mathbf{Y})$. In contrast, regression coefficients and other properties of regression models are generally not symmetric in the roles of the variables: $B_{\mathbf{Y}|\mathbf{X}} \neq B_{\mathbf{X}|\mathbf{Y}}$.

Table 3.2 Tests for Multivariate Hypotheses
 $\hat{\eta}_m =$ Strength of Multivariate Association of Test m
 $S_h = \hat{\Delta}$ and $S_e = (N - r)\hat{\Sigma}_*$

Name	Statistic	Principle	$\hat{\eta}_m$	Univariate Case
HLT	$\sum_{k=1}^s \frac{\hat{\rho}_k^2}{(1-\hat{\rho}_k^2)} = \text{tr}(S_h S_e^{-1})$	ANOVA analog	$\frac{\text{HLT}/s}{1 + \text{HLT}/s}$	$\frac{\hat{\rho}^2}{(1-\hat{\rho}^2)} = \frac{SSH}{SSE}$
PBT	$\sum_{k=1}^s \hat{\rho}_k^2 = \text{tr}[S_h(S_h + S_e)^{-1}]$	Substitution	$\frac{\text{PBT}}{s}$	$\hat{\rho}^2 = \frac{SSH}{SSH+SSE}$
WLK	$\prod_{k=1}^s (1-\hat{\rho}_k^2) = S_e(S_h + S_e)^{-1} $	Likelihood ratio	$1 - \text{WLK}^{1/g}$	$1 - \hat{\rho}^2 = \frac{SSE}{SSH+SSE}$
RLR	$\max_k \hat{\rho}_k^2 = \frac{\text{max eigenvalue}}{S_h(S_h + S_e)^{-1}}$	Union-intersection	$\hat{\rho}_1^2$	$\hat{\rho}^2 = \frac{SSH}{SSH+SSE}$
UNIREP	$\frac{\text{tr}(S_h)}{\text{tr}(S_e)}$	Best with sphericity	$\frac{\text{tr}(S_h)}{\text{tr}(S_h + S_e)}$	$\hat{\rho}^2 = \frac{SSH}{SSH+SSE}$

For consistency with the approximate distribution for WLK discussed in the remainder of the section,

$$g = \begin{cases} 1 & a^2 b^2 \leq 4 \\ [(a^2 b^2 - 4)/(a^2 + b^2 - 5)]^{1/2} & \text{otherwise.} \end{cases} \tag{3.39}$$

However, choosing $g = s$ leads to a simpler interpretation of $\hat{\eta}_{\text{WLK}}$, the measure of multivariate association, in terms of a geometric mean of (canonical) error variances. In fact, if $s \leq 2$, then $g = s$.

The $b \times b$ matrix $\hat{\Delta}$ contains the sums of squares for the hypothesis and reduces to the scalar sum-of-squares hypothesis, SSH , whenever $b = 1$. The $b \times b$ matrix $(N - r)\hat{\Sigma}_*$ contains the sums of squares for error and reduces to the scalar sum of squares SSE whenever $b = 1$. Also, $\hat{\Delta}[\hat{\Delta} + (N - r)\hat{\Sigma}_*]^{-1}$ reduces to the scalar $SSH/(SSH + SSE) = \hat{\rho}^2$, the estimated squared multiple correlation (which may or may not be adjusted for an intercept). Of course, any 1×1 matrix has only one eigenvalue, the scalar itself. Under the null, with $s = \min(a, b) = 1$, each of the four MULTIREP statistics can be expressed exactly as a one-to-one function of each other, and of an F random variable with numerator degrees of freedom $\nu_1 = a$ and denominator degrees of freedom $\nu_2 = N - r = \nu_e$.

Under the null, with $s > 1$, the MULTIREP statistics are not one-to-one functions of each other, and exact distributions are known only for special cases.

However, approximations based on an F random variable which match two moments are available for three of the four statistics, as detailed in Table 3.3. When $s > 1$ and $m \in \{\text{HLT}, \text{PBT}, \text{WLK}\}$, an approximate p value may be computed as follows. Denominator degrees of freedom (df) $\nu_2(m)$ are in Table 3.3. With $\nu_e = N - r$, the numerator degrees of freedom $\nu_1(m)$ are $\nu_1(\text{HLT}) = ab$, $\nu_1(\text{WLK}) = ab$, and

$$\nu_1(\text{PBT}) = ab \frac{1}{s(\nu_e + a)} \left[\frac{s(\nu_e + s - b)(\nu_e + a + 2)(\nu_e + a - 1)}{\nu_e(\nu_e + a - b)} - 2 \right]. \quad (3.40)$$

Table 3.3 Denominator df for F Approximations

Test	$\nu_2(m)$	Author
HLT	$\frac{[\nu_e^2 - \nu_e(2b + 3) + b(b + 3)](ab + 2)}{\nu_e(a + b + 1) - (a + 2b + b^2 - 1)} + 4$	McKeon (1974)
PBT	$\frac{\nu_e + s - b}{\nu_e + a} \left[\frac{s(\nu_e + s - b)(\nu_e + a + 2)(\nu_e + a - 1)}{\nu_e(\nu_e + a - b)} - 2 \right]$	Muller (1998)
WLK	$g[\nu_e - (b - a + 1)/2] - (ab - 2)/2$	Rao (1951)

Computing

$$f_{\text{obs}}(m) = \frac{\hat{\eta}_m / \nu_1(m)}{(1 - \hat{\eta}_m) / \nu_2(m)} \quad (3.41)$$

leads to the associated approximate p value

$$p(m) = 1 - F_F[f_{\text{obs}}(m); \nu_1(m), \nu_2(m)]. \quad (3.42)$$

Harris (1975, Appendix B) provided a useful method for directly approximating tail probabilities of RLR.

Eigenvalue k of $S_h S_e^{-1} = \hat{\Omega} / \nu_e$ is a one-to-one function of eigenvalue k of $S_h(S_h + S_e)^{-1}$ and also of eigenvalue k of $S_e(S_h + S_e)^{-1}$. Computational accuracy considerations lead to preferring to compute the eigenvalues of $S_h S_e^{-1}$, namely $\{\hat{\rho}_k^2 / (1 - \hat{\rho}_k^2)\}$. A standard and simple approach allows converting the nonsymmetric matrix to a symmetric matrix which has the same eigenvalues (and different eigenvectors which allow computing the eigenvectors of the original matrix). The Cholesky method, among other methods, allows finding $\hat{\Phi}$ such that $\hat{\Sigma}_* = \hat{\Phi} \hat{\Phi}'$, which implies $\hat{\Sigma}_*^{-1} = \hat{\Phi}^{-t} \hat{\Phi}^{-1}$. In turn, it is straightforward to prove that the eigenvalues of $\hat{\Omega}$ coincide with the eigenvalues of the symmetric matrix

$$\hat{\Omega}_\Phi = \hat{\Phi}^{-1} \hat{\Delta} \hat{\Phi}^{-t}. \quad (3.43)$$

Similarly, if $S_e = \nu_e \hat{\Sigma}_* = F_e' F_e'$, then the eigenvalues of $S_h S_e^{-1}$, namely $\{b_k\} = \{\hat{\rho}_k^2 / (1 - \hat{\rho}_k^2)\}$, coincide with the eigenvalues of

$$\widehat{\Omega}_\Phi/\nu_e = \mathbf{F}_e^{-1} \mathbf{S}_h \mathbf{F}_e^{-t}. \quad (3.44)$$

Obviously $\widehat{\rho}_k^2 = b_k/(1 + b_k)$ gives the squared canonical correlations, from which all of the test statistics may be computed, as detailed above.

3.9 COMPUTING UNIREP TESTS

The four UNIREP tests are *not* functions of the eigenvalues of $\mathbf{S}_h \mathbf{S}_e^{-1} = \widehat{\Delta}(\widehat{\Delta} + \nu_e \widehat{\Sigma}_*)^{-1}$. They all use the same test statistic,

$$f_{\text{obs}}(\mathbf{U}) = \frac{\text{tr}(\mathbf{S}_h)/(ab)}{\text{tr}(\mathbf{S}_e)/(b\nu_e)} = \frac{\text{tr}(\widehat{\Delta})/a}{\text{tr}(\widehat{\Sigma}_*)}, \quad (3.45)$$

and corresponding measure of multivariate association,

$$\widehat{\eta}_U = \frac{\text{tr}(\mathbf{S}_h)/\text{tr}(\mathbf{S}_e)}{1 + \text{tr}(\mathbf{S}_h)/\text{tr}(\mathbf{S}_e)} = \frac{\text{tr}(\mathbf{S}_h)}{\text{tr}(\mathbf{S}_h + \mathbf{S}_e)}. \quad (3.46)$$

The uncorrected test uses the p value

$$p(\text{Un}) = 1 - F_F[f_{\text{obs}}(\mathbf{U}); ab, b\nu_e]. \quad (3.47)$$

The Geisser-Greenhouse test reduces degrees of freedom by the maximum likelihood estimator of ϵ , namely $\widehat{\epsilon} = b^{-1} \text{tr}^2(\widehat{\Sigma}_*)/\text{tr}(\widehat{\Sigma}_*)^2$:

$$p(\text{GG}) = 1 - F_F[f_{\text{obs}}(\mathbf{U}); ab\widehat{\epsilon}, b\nu_e\widehat{\epsilon}]. \quad (3.48)$$

In seeking an approximately unbiased estimator, the Huynh-Feldt test uses $\widetilde{\epsilon} = (Nb\widehat{\epsilon} - 2)/[b(\nu_e - b\widehat{\epsilon})]$:

$$p(\text{HF}) = 1 - F_F[f_{\text{obs}}(\mathbf{U}); ab\widetilde{\epsilon}, b\nu_e\widetilde{\epsilon}]. \quad (3.49)$$

The fact that $\widetilde{\epsilon}$ may exceed 1.0 leads to using $\widetilde{\epsilon}_t = \min(\widetilde{\epsilon}, 1)$. The Box conservative test uses the lower bound for ϵ , namely $1/b$:

$$p(\text{Box}) = 1 - F_F[f_{\text{obs}}(\mathbf{U}); a, \nu_e]. \quad (3.50)$$

For data analysis, UNIREP tests differ only due to the degrees of freedom multipliers, which are always in the same order: Box, GG, HF, and uncorrected, with values $1/b \leq \widehat{\epsilon} \leq \widetilde{\epsilon}_t \leq 1$. Furthermore, the p values will always be in the reverse order.

If all $\lambda_k = \lambda_1$, then $\epsilon = 1$ and $\text{Dg}(\boldsymbol{\lambda}) = \lambda_1 \mathbf{I}_b$, which corresponds to a spherical Gaussian distribution. Under sphericity $f_{\text{obs}}(\mathbf{U}) \sim F\{ab, b\nu_e, \text{tr}(\boldsymbol{\Omega})\}$ (exactly), the test is exactly size α and uniformly most powerful (among similarly invariant tests). Box (1954a, b) observed that $1/b \leq \epsilon \leq 1$ and that $\epsilon < 1$ implies, under the null, $f_{\text{obs}}(\mathbf{U}) \sim F(ab\epsilon, b\nu_e\epsilon)$, an *approximate* result.

Muller, LaVange, Ramey, and Ramey (1992) reviewed power approximation for both UNIREP and MULTIREP tests. We leave the topic until Chapter 21.

3.10 CONFIDENCE REGIONS FOR Θ

As noted in Section 2.10, confidence regions can be obtained by inverting hypothesis tests, and a confidence region can be inverted to yield a hypothesis test. The definition in Section 2.10 can be extended easily from vector data \mathbf{y} and vector parameter θ (or σ^2) to matrix data, \mathbf{Y} , and matrix parameter Θ (or Σ). Doing so merely requires replacing \mathbf{y} by $\text{vec}(\mathbf{Y})$ and θ by $\text{vec}(\Theta)$ [or σ^2 by $\text{vec}(\Sigma)$] in Definition 2.12. In Section 15.6 we prove a variety of results for confidence intervals in univariate models. In Section 16.10 we describe extensions to multivariate models.

3.11 SUFFICIENT STATISTICS FOR THE MULTIVARIATE MODEL

If the multivariate linear model $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with Gaussian distribution is correct, then the matrix

$$\begin{aligned} \mathbf{S} &= \begin{bmatrix} \mathbf{X}'\mathbf{X} & \mathbf{X}'\mathbf{Y} \\ \mathbf{Y}'\mathbf{X} & \mathbf{Y}'\mathbf{Y} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{X}' \\ \mathbf{Y}' \end{bmatrix} [\mathbf{X} \ \mathbf{Y}] \end{aligned} \quad (3.51)$$

contains all of the complete sufficient statistics. If \mathbf{X} contains an intercept, then, without loss of generality, \mathbf{X} may be arranged with the intercept in column 1, with $\mathbf{X} = [\mathbf{1}_N \ \mathbf{X}_1]$, for \mathbf{X}_1 of dimension $N \times (q - 1)$. In turn,

$$\begin{aligned} \mathbf{S} &= \begin{bmatrix} \mathbf{1}'\mathbf{1} & \mathbf{1}'\mathbf{X}_1 & \mathbf{1}'\mathbf{Y} \\ \mathbf{X}_1'\mathbf{1} & \mathbf{X}_1'\mathbf{X}_1 & \mathbf{X}_1'\mathbf{Y} \\ \mathbf{Y}'\mathbf{1} & \mathbf{Y}'\mathbf{X}_1 & \mathbf{Y}'\mathbf{Y} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{1}' \\ \mathbf{X}_1' \\ \mathbf{Y}' \end{bmatrix} [\mathbf{1} \ \mathbf{X}_1 \ \mathbf{Y}] \end{aligned} \quad (3.52)$$

contains all of the complete sufficient statistics for estimation of all parameters of the regression models (one for each column of \mathbf{Y}) identified by the relationship

$$E(\mathbf{Y}|\mathbf{X}). \quad (3.53)$$

Here $[\mathbf{1} \ \mathbf{X}_1 \ \mathbf{Y}]$ is $N \times (q + p)$ so \mathbf{S} is $(q + p) \times (q + p)$. All parameter estimators and general linear hypothesis tests (both MULTIREP and UNIREP) depend on the data only through the elements of \mathbf{S} . Conveniently, the raw data are not needed for parameter estimation or testing the general linear hypothesis.

3.12 ALLOWING MISSING DATA IN THE MULTIVARIATE MODEL

Definition 3.9 (a) If $(1 \times p_i)$ random matrix $\mathbf{Y}_i = \{y_{ij}\}$ represents the potential response of ISU i , then an element y_{ij} is said to be *missing* if a realized value for y_{ij} is not included in the statistical analysis.

(b) If only two patterns of data occur, namely either all elements of \mathbf{Y}_i are present or all elements are missing, the observations are described as having only *casewise missing* observations.

Examples of missing values include outlier values intentionally omitted by the analyst, recorded values that were lost during data entry, values that were never recorded because the ISU (individual participant) was not available for evaluation, and interval-censored values treated as unknown by the analyst. The definition applies to the multivariate model ($p_i \equiv p \forall i$) as well as to clustered data in general.

The standard methods described in the present chapter conveniently allow casewise missing observations in the multivariate (and univariate) linear model. If the mechanism causing data to be missing does not lead to selection biases, then the approach gives optimal estimators and exact tests. Partially missing \mathbf{X}_i or \mathbf{Y}_i do not have such nice properties. In the multivariate linear model, nearly all research on the topic has focused on having complete \mathbf{X} and partially missing \mathbf{Y}_i . Such patterns occur naturally in an experiment with random assignment to treatment, as in a typical clinical trial. Little (1992) reviewed methods for regression with incomplete \mathbf{X} . Little and Rubin (2002) reviewed methods for missing data, with particular emphasis on estimation from a likelihood perspective.

Using the notation of Stewart (2000) and others, a binary random variable $r_{ij} = 1$ indicates y_{ij} is not missing while $r_{ij} = 0$ indicates that it is missing. With $(1 \times p_i)$ matrix \mathbf{R}_i containing all r_{ij} for ISU i , specifying the entire response for ISU i requires knowing $\{\mathbf{Y}_i, \mathbf{R}_i\}$. Vertical concatenation of $\{\mathbf{Y}'_i\}$ and similar concatenation of $\{\mathbf{R}'_i\}$ yield the pair of random vectors $\{\mathbf{y}, \mathbf{r}\}$. It is assumed that an appropriate parametric statistical model can be formulated for \mathbf{y} . If \mathbf{y} has a density function, then the model can be represented as $\{\mathbf{y}, f_{\mathbf{y}}(\mathbf{y}_*|\mathbf{X}; \boldsymbol{\theta}), \boldsymbol{\theta} \in \Theta\}$.

Information for estimation and inferences about $\boldsymbol{\theta}$ is available only via $\{\mathbf{y}_{\text{obs}}, \mathbf{R}\}$, in which \mathbf{y}_{obs} is the vector of nonmissing values. The missing values are denoted \mathbf{y}_{mis} . Usually, most or all of this information about $\boldsymbol{\theta}$ is anticipated to come from \mathbf{y}_{obs} . If the observed pattern of missing values \mathbf{R}_* contains no additional information about $\boldsymbol{\theta}$, then the underlying missing-data mechanism is said to be ignorable as defined by Rubin (1976). We then say the missing values are ignorably missing.

In general, the likelihood function (any function proportional to the joint density function of $\{\mathbf{y}_{\text{obs}}, \mathbf{r}\}$) is obtained by an integration over possible values of \mathbf{y}_{mis}

(Little and Rubin, 2002). The process creates a marginal density:

$$f_{\mathbf{y}}(\mathbf{y}_{*obs}, \mathbf{r}_{*} | \mathbf{X}; \boldsymbol{\theta}, \boldsymbol{\psi}) = \int f_{\mathbf{y}}(\mathbf{y}_{*obs}, \mathbf{y}_{*mis} | \mathbf{X}; \boldsymbol{\theta}) f_{\mathbf{r}}(\mathbf{r}_{*} | \mathbf{y}_{*obs}, \mathbf{X}; \boldsymbol{\psi}, \boldsymbol{\theta}) d\mathbf{y}_{*mis}. \quad (3.54)$$

Definition 3.10 The following definitions follow Rubin (1976).

(a) If \mathbf{r} and \mathbf{y} are statistically independent, given \mathbf{X} , then the joint density is

$$f_{\mathbf{y}}(\mathbf{y}_{*obs}, \mathbf{r}_{*} | \mathbf{X}; \boldsymbol{\theta}, \boldsymbol{\psi}) = f_{\mathbf{y}}(\mathbf{y}_{*obs} | \mathbf{X}; \boldsymbol{\theta}) f_{\mathbf{r}}(\mathbf{r}_{*} | \mathbf{X}; \boldsymbol{\psi}, \boldsymbol{\theta}), \quad (3.55)$$

and the missing values are said to be *missing completely at random* (MCAR).

(b) The missing-data mechanism is said to be *ignorable* if the following factorization holds true:

$$f_{\mathbf{y}}(\mathbf{y}_{*obs}, \mathbf{r}_{*} | \mathbf{X}; \boldsymbol{\theta}, \boldsymbol{\psi}) = f_{\mathbf{y}}(\mathbf{y}_{*obs} | \mathbf{X}; \boldsymbol{\theta}) f_{\mathbf{r}}(\mathbf{r}_{*} | \mathbf{y}_{*obs}, \mathbf{X}; \boldsymbol{\psi}). \quad (3.56)$$

(c) Invariance of $f_{\mathbf{r}}(\mathbf{r}_{*} | \mathbf{y}_{*obs}, \mathbf{y}_{*mis}, \mathbf{X}; \boldsymbol{\psi}, \boldsymbol{\theta})$ to the possible values of \mathbf{y}_{*mis} gives data described as *missing at random* (MAR).

The definitions lead to the following observations. If the conditional distribution of \mathbf{r} in equation (3.55) does not depend on $\boldsymbol{\theta}$, then no information about $\boldsymbol{\theta}$ is neglected if \mathbf{r} is ignored. Rubin (1976) proved that (3.56) holds if and only if $f_{\mathbf{r}}(\mathbf{r}_{*} | \mathbf{y}_{*obs}, \mathbf{y}_{*mis}, \mathbf{X}; \boldsymbol{\psi}, \boldsymbol{\theta})$ is invariant to the possible values of $\{\mathbf{y}_{*mis}, \boldsymbol{\theta}\}$ when evaluated at $\{\mathbf{r}_{*}, \mathbf{y}_{*obs}\}$. The invariance with respect to $\boldsymbol{\theta}$ is a condition Rubin (1976) referred to as the absence of a priori ties between the parameters of the two densities in the factorization.

If factorization (3.55) or (3.56) holds, then maximizing the $\{\mathbf{y}_{obs}, \mathbf{r}\}$ likelihood with respect to $\boldsymbol{\theta}$ is equivalent to maximizing $f_{\mathbf{y}}(\mathbf{y}_{*obs} | \mathbf{X}; \boldsymbol{\theta})$. It is in this sense that \mathbf{r} can be ignored. It is important to realize, however, that maximization of a likelihood can yield parameter estimates and likelihood ratios without providing standard errors for the parameter estimators. Verbeke and Molenberghs (2000, Chapter 21) gave examples of bias stemming from using the expected information matrix for approximating standard errors in the MAR case. Heitjan (1994) gave a clear statement of the problem. Diggle and Kenward (1994), Stewart (2000), and Lipsitz et al. (2002) discussed illustrative MAR examples.

With all assumptions of the $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors met, having MAR or MCAR data allows computing maximum likelihood estimates of $\{\mathbf{B}, \boldsymbol{\Sigma}\}$, although usually through iterative methods. One convenient approach uses the estimates in complete data formulas for confidence intervals and testing hypotheses. Unfortunately, Barton and Cramer (1989) demonstrated that the naive approach leads to optimistic (biased) estimates of precision with small to moderate sample sizes and MCAR data. The same authors, as well as Catellier and Muller (2000), described very simple adjustments to the degrees of freedom for the MULTIREP and UNIREP tests in the MCAR case. In simulations, the adjusted tests completely or nearly control test size, even in very small samples ($N = 12$

and $p = 6$). In contrast, standard linear mixed model tests greatly inflated test size (up to 0.40 with a nominal value of $\alpha = 0.05$). The Appendix (Section A.2) contains descriptions of free SAS/IML[®] code, and where to find it on the Web, which implements the methods.

Informatively missing data present one of the most vexing problems in statistical modeling. Almost unavoidably, quite specific assumptions must be made which derive from a particular scientific setting and analysis goal. As noted earlier, Little and Rubin (2002) provided the best starting point for further reading. General techniques have not been developed, and new developments continue. An analyst seeking the best method available for a specific analysis would be wise to carefully review the statistical literature. Using the Current Index to Statistics and other electronic databases greatly eases the pain of the search.

EXERCISES

All exercises refer to the breast cancer example described in the Appendix (Section A.1). Use P0104.SD2 for any data analysis. When necessary, consider “Benign” as the reference cell, which corresponds to choosing “MALIGN” as a predictor. Use a nominal test size of 0.025 (chosen due to conducting two planned analyses in the original study).

When needing to understand a test or contrast, one can often ignore the multivariate nature of the design and apply the logic of univariate ANOVA design and interpretation. (The approach fails for derivations of distribution theory.) To specify a within-subject contrast matrix, it may help to first specify a contrast matrix for between-subject effects based on cell mean coding (with the correct dimensions and factor pattern) and then transpose it.

3.1 One focus of the study was testing the hypothesis of no difference in the variables DLOGROI1–DLOGROI3 DLOG_P_1–DLOG_P_3 between patients with malignant and benign pathology. Briefly specify an appropriate multivariate linear model $GLM_{N,q,p}(Y_i; X_i B, \Sigma)$ with Gaussian errors. Use reference cell coding. Include values of all dimensions for the specific data in hand as well as brief definitions of parameter matrix elements.

3.2.1 For the model chosen in exercise 3.1, explicitly specify all contrast matrices needed to test a “MANOVA” hypothesis of no differences between benign and malignant. Include all dimensions for the specific data in hand and brief definitions of parameter matrix elements.

3.2.2 Choose a test statistic for 3.2.1 and briefly justify your choice. This must be done a priori (before looking at *any* data).

3.2.3 Assuming that the overall MANOVA test just discussed is significant, explicitly specify a modest number of scientifically interesting and appropriate stepdown tests and associated contrast matrices. It is acceptable to use the simplest approach to control multiple testing bias, namely a Bonferroni correction. Include all dimensions for the specific data in hand and brief definitions of parameter matrix elements.

3.3 Treating tissue type [(region of interest (ROI), versus parenchyma] and time (with three levels) as factors in a factorial design with benign versus malignant may have more appeal than a MANOVA analysis.

3.3.1 Use cell mean style coding in this exercise. Specify an appropriate multivariate model.

3.3.2 Assuming a factorial approach describe an appropriate “source table” by listing the sources to be tested, and associated hypothesis degrees of freedom.

3.3.3 Explicitly specify all contrast matrices for each source listed in 3.1. Include all dimensions for the specific data in hand and brief definitions of parameter matrix elements. Use individual rows (or columns) that test *differences in means* (which are the secondary parameters).

3.4 Choose a test statistic for 3.1 and 3.2 and briefly justify your choice. This must be done a priori (before looking at *any* data).

3.5 Explicitly specify all contrast matrices for each source listed in 3.1. Include all dimensions for the specific data in hand and brief definitions of parameter matrix elements. Use individual rows (or columns) that test *polynomial trends* which are the secondary parameters, even though not all factors have levels defined in terms of a continuous variable.

3.6 Using PROC GLM and the data supplied, implement the two-way model described in 3.1, 3.3, and 3.4. *Hint*: Use a particular statement type available with GLM which does most of the coding and testing work associated with repeated measures automatically.

3.6.1 Provide sufficient source code and a compact numerical version of a source table. Include numerator degrees of freedom, a test statistic p value, and an appropriate measure of multivariate association.

3.6.2 Provide a brief scientific interpretation of the results.

3.7 Construct three new difference variables at times 1, 2, and 3: DLOGROI1–DLOG_P_1, DLOGROI2–DLOG_P_2, and DLOGROI3–DLOG_P_3.

3.7.1 Fit a one-way model, again using the same statement approach in PROC GLM, with Time as a factor. Report and interpret an appropriate test of no difference in the constructed variables between patients with malignant and benign pathology.

3.7.2 How does the model relate to the two-way model in exercise 3.4?

3.8 Examine the maximum likelihood estimate of ϵ for the UNIREP tests. Compare the p values of the four UNIREP tests and also the four MULTIREP tests. Make a recommendation for a choice of statistic in future studies of the same sort.

3.9 (*optional, noncredit*) Use LINMOD to repeat the analysis. Select appropriate options to enrich the output and help understanding. If you used LINMOD in the first place, then use PROC GLM or a procedure in another computer language. You should be able to reproduce almost exactly nearly all values [except for some multivariate p values if $s = \min(a, b) > 1$].

The Appendix (Section A.2) contains a brief description of the free software LINMOD and where to find it on the Web.

CHAPTER 4

Generalizations of the Multivariate Linear Model

4.1 MOTIVATION

As mentioned in the previous chapter, the multivariate general linear model has a number of limitations: The multivariate model does not directly tolerate incomplete or mistimed data; the multivariate model does not allow the design matrix to vary across responses; the multivariate model does not explicitly allow modeling the covariance structure. A number of generalizations of the multivariate model have been developed to avoid the limitations.

In the present chapter we briefly survey some generalizations of the multivariate linear model and its special case, the univariate linear model. Many of them stand somewhere in between the mixed and multivariate linear models, in terms of both theory and applications. As a broad generality, all provide well-behaved estimation but may have difficulty providing completely accurate inference in small samples. However, the alternative use of a mixed model may provide even less accuracy, depending on the situation. Lacking exact and perfect methods, good statistical practice, as always, centers on using the best available approximation.

In contrast, the general linear mixed model has none of the limitations. Unfortunately, the generality of the mixed model may come at a steep price. Even with a correctly specified model, the approach can lead to extremely inaccurate inference (optimistically small confidence intervals and inflated test size), especially with small to moderate sample sizes. The inaccuracy arises from what Littel (2003) described as “approximations piled on approximations.” Furthermore, limitations of current methods make it difficult to check the validity of the model, especially the covariance component. Simulation results make it clear that misspecification of the covariance model may introduce (additional) substantial inaccuracy in inference (Park, Park, and Davis, 2001; Muller, Edwards, Simpson, and Taylor, 2006).

Although linear in the expected-value parameters, the likelihood varies nonlinearly as a function of the covariance parameters. As a consequence, computing estimates for a linear mixed model requires iterative solution of a system of simultaneous nonlinear equations. Collinearity arising from less than

careful attention to data scaling or location and less-than-full-rank coding schemes easily and often degrades the performance of the algorithms needed. The resulting difficulties in computing estimates can lead the unsophisticated user to propose a valid model which fails to converge. Simplifying the covariance model greatly increases the chances of achieving convergence. However, the simplification may also misspecify the covariance model and thereby introduce severe bias.

A simple example illustrates the concern. Clinical trials of a pharmaceutical agent which involve repeated measures frequently lead to some missing and mistimed data, which make computing estimates for linear mixed models more difficult. In many such cases, analysts have assumed compound symmetry of covariance. Although extremely convenient (due to helping convergence), the assumption often seems implausible for a sequence of (time) ordered responses.

4.2 THE GENERALIZED GENERAL LINEAR UNIVARIATE MODEL: EXACT AND APPROXIMATE WEIGHTED LEAST SQUARES

The $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ can be generalized in many ways. The motivation lies in the need to allow for patterns of dependence, rather than complete independence, among response values. The simplest way to allow such dependence is to assume $\mathcal{V}(\mathbf{y}) = \mathcal{V}(\mathbf{e}) = \boldsymbol{\Upsilon}$, with $\boldsymbol{\Upsilon}$ of dimension $N \times N$, symmetric, and positive definite or positive semidefinite, (which allows any covariance matrix). In contrast, however, the approach allows describing only extremely limited results of little practical value. The limitations arise from the fact that $\boldsymbol{\Upsilon}$ has $N(N+1)/2$ distinct parameters, which exceeds N , the number of observations. Increasing sample size only worsens the problem because the number of parameters increases more rapidly than N . Even assuming complete independence while allowing complete heterogeneity does not solve the problem. In that case $\boldsymbol{\Upsilon} = Dg(\mathbf{v})$ has N parameters, which implies increasing sample size never allows the number of observations to exceed the number of parameters to be estimated. Avoiding the problem requires adding assumptions which impose structure on $\boldsymbol{\Upsilon}$. Necessarily the number of parameters must grow more slowly than the sample size if reasonable estimators are to exist. The multivariate and mixed linear models both generalize the univariate linear in the same fashion, although in very different directions. Subsequent chapters are devoted to properties of two approaches.

Following McCullagh and Nelder (1989), the term “generalized linear model” refers to a model with expected values linear or loglinear in the parameters and the response distribution any member of the exponential family, not just the Gaussian. To avoid confusion, we introduce the following definitions.

Definition 4.1 A *generalized GLM* is indicated by the notation $GGLM_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \boldsymbol{\alpha}, \boldsymbol{\Upsilon})$, which describes all observations, not just a single observation for an independent sampling unit. The assumptions differ from a $GLM_{N,q}()$ in only one way, which is important. The assumption

“elements of the $N \times 1$ random vector \mathbf{y} are mutually independent.” is replaced by the assumption “elements of the $N \times 1$ random vector \mathbf{y} have (finite) $N \times N$, constant, covariance matrix Υ , which may be at least partially known.”

The $\text{GGLM}_{N,q}()$ notation describes all N observations at once because they may not be independent, while all other model notations used here describe the observations for a single independent sampling unit. In fact, if all off-diagonal elements of a full-rank Υ are nonzero, then a GGLM has only one independent sampling unit. In the special case of a grand mean model, which has $\mathbf{X}\beta = \mathbf{1}_N\mu$, the $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{1}_N\mu, \Upsilon)$ in many ways corresponds to a multivariate model with one observation, $\text{GLM}_{1,N,1}(\mathbf{y}'; \mathbf{1}'_N\mu, \Upsilon)$. With Gaussian data, the special case has $\mathbf{y} \sim \mathcal{N}_N(\mathbf{1}_N, \Upsilon)$.

Definition 4.2 A GGLM with Gaussian errors refers to a setting in which $\mathbf{y} \sim \mathcal{N}_N(\mathbf{X}\beta | \mathbf{R}\beta = \mathbf{a}, \Upsilon)$, which is an assumption of joint (“multivariate”) Gaussian distribution, not merely marginally Gaussian $\{y_i\}$. As for a univariate GLM, a GGLM may be described as either FR or LTFR and restricted or unrestricted, depending on \mathbf{X} , \mathbf{R} , and \mathbf{a} .

Lemma 4.1 For $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\beta | \mathbf{R}\beta = \mathbf{a}, \sigma^2 \mathbf{I}_N)$, elements of $\Upsilon = \mathcal{V}(\mathbf{y}) = \mathcal{V}(\mathbf{e})$ control many properties of the model.

- (a) If any off-diagonal element of Υ is not equal to zero, the model does not have independent observations.
- (b) If Υ has two or more distinct diagonal elements, the model does not have homogeneity of variance.
- (c) If $\Upsilon = \sigma^2 \mathbf{I}_N$, then the model meets the assumptions of the univariate $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta | \mathbf{R}\beta = \mathbf{a}, \sigma^2)$, which means the univariate $\text{GLM}_{N,q}()$ is a special case, namely $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\beta | \mathbf{R}\beta = \mathbf{a}, \sigma^2 \mathbf{I}_N)$.
- (d) More generally, the stringent condition $\Upsilon = \sigma^2 \mathbf{D}$ with \mathbf{D} symmetric, positive definite or semidefinite of rank $N_1 \leq N$, *known*, and not needing to be estimated allows converting a $\text{GGLM}_{N,q}()$ with Gaussian errors to a univariate $\text{GLM}_{N_1,q}()$ with Gaussian errors.

Proof. Parts (a), (b), and (c) are true due to properties of second moments. Part (d) has $0 < \text{rank}(\mathbf{D}) = N_1 \leq N$. With $\mathbf{V}'_1\mathbf{V}_1 = \mathbf{I}_{N_1}$, spectral decomposition gives $\mathbf{D} = \mathbf{V}_1 \text{Dg}(\mathbf{d}_1) \mathbf{V}'_1$. If $\mathbf{F}' = \text{Dg}(\mathbf{d}_1)^{-1/2} \mathbf{V}'_1$, then $\mathbf{D}^+ = \mathbf{V}_1 \text{Dg}(\mathbf{d}_1)^{-1} \mathbf{V}'_1 = \mathbf{F}\mathbf{F}'$. The original data satisfy the equation $\mathbf{y} = \mathbf{X}\beta + \mathbf{e}$, with $\mathbf{e} \sim (\mathcal{S})\mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{D})$. In turn, knowing \mathbf{D} allows transforming the model:

$$\mathbf{F}'\mathbf{y} = \mathbf{F}'\mathbf{X}\beta + \mathbf{F}'\mathbf{e} \tag{4.1}$$

$$\mathbf{y}_F = \mathbf{X}_F\beta + \mathbf{e}_F. \tag{4.2}$$

The relationship $\mathbf{F}'\mathbf{D}\mathbf{F} = \mathbf{I}_{N_1}$ allows concluding that $\mathbf{e}_F \sim \mathcal{N}_{N_1}(\mathbf{0}, \sigma^2 \mathbf{I}_{N_1})$. Any jointly Gaussian variables with zero covariance are statistically independent.

Therefore the $\{y_{Fi}\}$ meet the assumptions of $\text{GLM}_{N,q}(y_{Fi}; \mathbf{X}_{Fi}\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$ with Gaussian errors. \square

Definition 4.3 (a) Part (d) of the last lemma and the associated proof defines a model and data analysis process often described as *exact weighted least squares* or just *weighted least squares (WLS)*.

(b) Using observed data to estimate any property of \mathbf{D} gives *approximate weighted least squares (AWLS)* analysis.

(c) *Iterated approximate least squares (ITAWLS)* centers on alternately updating the estimates of mean and covariance parameters.

Such methods apply whenever all assumptions of a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ hold except homogeneity of variance, and $\mathcal{V}(y_i | \mathbf{X}_i\boldsymbol{\beta}) = \sigma_i^2 = \sigma^2 w_i$, with $\{w_i\}$ known and σ^2 unknown. As discussed in later chapters, the model can be transformed exactly to a GLM with all assumptions holding, including homogeneity. In such cases the transformation leads to optimal estimators and exact (and optimal) tests with respect to the parameters of the original model.

To be precise, weighted least squares may be referred to as exact weighted least squares in order to distinguish it from approximate weighted least squares (AWLS), which estimates some features of \mathbf{D} in estimating $\boldsymbol{\Upsilon} = \sigma^2 \mathbf{D}$. Iterated approximate least squares (ITAWLS) relies on an alternating updating of the estimates of mean and covariance parameters. With Gaussian errors, in many settings the approach leads to maximum likelihood estimation. Not surprisingly, without appropriate adjustments, inference may be inaccurate in small to moderate sample sizes.

Many popular statistical methods, including most mixed model analyses, implicitly include some variation of ITAWLS or AWLS and then operate as though the estimated covariance structure was the population structure. In small to moderate sample sizes the approach can lead to substantial optimistic bias in confidence intervals and hypothesis tests. We urge the reader to always distinguish between the two approaches in reading and evaluating the work of others and in reporting analyses using either method. Including some indication of the expected impact of estimating the weights seems necessary for nonstatisticians to appreciate the amount of uncertainty introduced by an approximate analysis.

Example 4.1 Most methods for confidence intervals and hypothesis tests in current mixed model software implicitly depend on a large-sample assumption: Using covariance estimates in weighted least squares forms introduces no bias. Later chapters centered on mixed models include the details.

4.3 DOUBLY MULTIVARIATE MODELS

The standard multivariate linear applies to the doubly multivariate setting and provides exact size- α tests. However, with p response variables, each measured at t times, the approach implies estimating a $(pt) \times (pt)$ fully unstructured covariance matrix with $pt(pt + 1)/2$ distinct elements. The structure of the data makes it very appealing to assume the covariance matrix equals $\Sigma_1 \otimes \Sigma_2$, with Σ_1 the $p \times p$ covariance among responses and Σ_2 the $t \times t$ covariance among times. The direct-product form has far fewer covariance parameters, $p(p + 1)/2 + t(t + 1)/2$. Even in the simplest case with $p = t = 2$, the unstructured model has 10 covariance parameters, while the direct-product form has 6. In turn, $p = t = 5$ gives 325 versus 111 covariance parameters. Timm (2002) reviewed doubly multivariate models based on the direct-product covariance assumption. The work of Boik deserves special attention.

4.4 SEEMINGLY UNRELATED REGRESSIONS

Definition 4.4 A *seemingly unrelated regression model* (sometimes called a multiple design matrix model) corresponds to a $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ generalized to allow the design matrix and associated parameters to vary across columns of Y_i .

Much of the work on seemingly unrelated regression was motivated by econometric applications. Srivastava and Giles (1987) provided a book-length treatment.

Repeated-measures settings may naturally lead to the desire to vary the design matrix across response values. A clinical trial of a drug in the elderly may need to use the dose of a second drug as a covariate. The dose of the second drug may represent a time-varying covariate (or repeated covariate). Two natural variations occur. In the first, only the contemporaneous dose of the second drug matters, while in the second the contemporaneous and all preceding doses matter. If cast as a $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$, the first setting implies the desire for a block of the regression coefficients matrix B to be diagonal. The second setting implies the desire for a block of B to be triangular. Either condition requires imposing nonlinear constraints on B [which correspond to linear constraints on $\text{vec}(B)$].

4.5 GROWTH CURVE MODELS (GMANOVA)

When growth is observed over time by repeated measurement of a characteristic, the recorded longitudinal pattern can be plotted in two dimensions as a response-versus-time growth curve. The ordered responses of interest might be childrens' linearly increasing weights recorded at ages 1, 2.5, 3, and 4.5 months.

Alternately, the ordered responses might be longitudinal measurements of a drug's exponentially decreasing serum concentration in healthy volunteers. By definition, growth curve model (GCM) analysis focuses on investigating the functional relationship among ordered responses. Conventional GCM methods apply to growth data (indexed by *time* or *age*) and to other analogs such as dose-response data (indexed by *dose*), location-response data (indexed by *distance*), or response-surface data (indexed by two or more variables such as *latitude* and *longitude*), for example. Growth data may exhibit either positive or negative growth, as in the case of the rise or decline of bacterial colonies grown in laboratory dishes. Although most applications of GCM methods center on longitudinal observations on a one-dimensional characteristic (e.g., weight of children), the methods can also apply to multidimensional characteristics such as {weight, height}. The GCM discussed in the present chapter is the classical model considered by Potthoff and Roy (1964), Grizzle and Allen (1969) and Rao (1973). It is also known as a GMANOVA model. Kshirsagar and Smith (1995) provided the best single source.

The scope of our discussion will focus on any such one-dimensional collection of ordered responses with *consistently timed observations*; that is, all the ISUs studied have been observed (or not) on the same occasions (*ages* or *times* or *doses*, etc.). In such cases the observational design must have specified recording the response of interest at p different *times*, $\{t_1, t_2, \dots, t_p\}$, *doses*, $\{d_1, d_2, \dots, d_p\}$, etc. The observation process creates a matrix of responses, \mathbf{Y} ($N \times p$). If $N = 1$ and the response of interest is a child's weight, then plotting weight at several ages indicates a temporal pattern of growth. A univariate linear model for weight given age could be fitted with a design matrix \mathbf{T} ($p \times m$) expressing the child's central tendency as a linear or curvilinear function of age. Here \mathbf{T} is an example of a *within-subject* design matrix. If $N > 1$, then a separate curve could be fitted for each child to obtain a separate ($m \times 1$) matrix of regression parameter estimators for each ISU, $\{\hat{\mathbf{B}}_i = \mathbf{Y}_i \mathbf{T} (\mathbf{T}' \mathbf{T})^{-1} : i \in \{1, 2, \dots, N\}\}$. A simple average of the N fitted curves is a proper (if not efficient) estimator of the population growth curve: $\hat{\mathbf{B}} = (\hat{\mathbf{B}}_1 + \hat{\mathbf{B}}_2 + \dots + \hat{\mathbf{B}}_N) / N$. In the following, \mathbf{X} ($N \times q$) represents a *between-subject* design matrix, which contains intersubject explanatory variables such as gender. The ($q \times m$) efficient estimator has the form

$$\hat{\mathbf{B}} = [(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{Y} [\mathbf{T} (\mathbf{T}' \mathbf{T})^{-1}]]. \quad (4.3)$$

If the subjects are a homogeneous group, then $\mathbf{X} = \mathbf{1}$ ($N \times 1$) is the appropriate choice for computing $\hat{\mathbf{B}}$. The choice of \mathbf{T} defines the functional form of the population growth curve by describing a functional relationship between weight and age. It also defines functional dependencies among the age-specific mean heights. If mean height is linear in age, then the mean at 3.5 years is constrained by linearity to lie halfway between the means at 2.5 and 4.5 years. Necessarily the means are collinear. The GCM explicitly addresses the dependencies via the within-subject design matrix \mathbf{T} and simultaneously addresses intersubject factors via the between-subject design matrix \mathbf{X} .

The GCM can be fitted using methods of either a restricted multivariate linear model for \mathbf{Y} ($N \times p$) or a linear mixed model for $\text{vec}(\mathbf{Y})$. The GCM thus provides common ground for our discussion of connections between multivariate linear models and linear mixed models.

For \mathbf{Y}_i ($1 \times p$) representing growth or dose-response in individual i , a preliminary multivariate model incorporating only a between-subject design matrix \mathbf{X} is $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}_{Y|X}, \Sigma_{Y|X})$. The key feature of the general model is that the $(1 \times p)$ matrix of means for individual i , $\mathbf{B}_i = \boldsymbol{\mu}'_i = E(\mathbf{Y}_i) = \mathbf{X}_i \mathbf{B}_{Y|X}$, is constrained to be a linear combination of the columns of \mathbf{X}_i ($1 \times r$). In turn, the GCM is only of interest when the variation of the central tendency within individual i as a function of time (dose) satisfies a linear model, $E(\mathbf{Y}'_i) = \mathbf{T}' \mathbf{B}_{Y|T}$. It follows that we also wish to constrain $\boldsymbol{\mu}'_i$ to be a linear combination of the rows of \mathbf{T} ($m \times p$), $\boldsymbol{\mu}'_i = E(\mathbf{Y}_i) = \mathbf{B}_{Y|T} \mathbf{T}$. The constraint can be imposed by modifying the preliminary multivariate model with the restrictions $m \leq p$ and

$$\mathbf{B}_{Y|X} \underset{q \times p}{=} \underset{(q \times m)(m \times p)}{\mathbf{B} \mathbf{T}} \tag{4.4}$$

With \mathbf{T} treated as a given constant, the constrained model may be represented as

$$\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}_{Y|X} | \mathbf{B}_{Y|X} = \mathbf{B} \mathbf{T}, \Sigma_{Y|X}). \tag{4.5}$$

We assign a special notation to the corresponding restricted multivariate GLM in the following definition. Hopefully $|\Sigma_{Y|X,T}|$ will be noticeably smaller than $|\Sigma_{Y|X}|$ because conditioning on both \mathbf{X} and \mathbf{T} should partially account for both inter-individual variance and intra-individual variance.

Definition 4.5 A *growth curve model* will be indicated by $\text{GCM}_{N,p,q,m}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B} \mathbf{T}, \Sigma)$ and includes the following assumptions.

1. The rows of the $N \times p$ random matrix \mathbf{Y} are mutually independent. With $\mathbf{Y}_i = \text{row}_i(\mathbf{Y}) = [y_{i1} \ y_{i2} \ \cdots \ y_{ip}]$, columns correspond to p ordered responses for p doses, times, etc., arrayed in vector $\mathbf{d} = [d_1 \ d_2 \ \cdots \ d_p]'$.
2. *Within-subject design matrix*, \mathbf{T} ($m \times p$), has $\text{rank}(\mathbf{T}) = m \leq p$ and is a fixed, known function of \mathbf{d} , and known without appreciable error. Consistent timing ensures \mathbf{d} and \mathbf{T} are constant $\forall i$.
3. With $\mathbf{X}_i = \text{row}_i(\mathbf{X})$, the $N \times q$ *between-subjects design matrix* \mathbf{X} has $\text{rank}(\mathbf{X}) = r \leq q \leq N$, and is fixed and known without appreciable error, conditional on knowing the sampling units, for data analysis.
4. Elements of \mathbf{B} ($q \times m$) are fixed and unknown regression coefficients.
5. The mean of \mathbf{Y}_i is $E(\mathbf{Y}_i | \mathbf{X}_i, \mathbf{T}) = \mathbf{X}_i \mathbf{B} \mathbf{T}$ ($1 \times q$)($q \times m$)($m \times p$).
6. Response \mathbf{Y}_i ($1 \times p$) has finite covariance matrix, Σ ($p \times p$), which is fixed, unknown, and positive definite or positive semidefinite. Also $\mathcal{V}[\text{vec}(\mathbf{Y}') | \mathbf{X}, \mathbf{T}] = \mathbf{I} \otimes \Sigma$.

Writing $\text{GCM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B} \mathbf{T} | \mathbf{R}_x \mathbf{B} \mathbf{R}_y = \mathbf{A}, \Sigma)$ specifies additional

explicit restrictions on parameters in \mathbf{B} through the fixed and known constants \mathbf{R}_x , \mathbf{R}_y , and \mathbf{A} .

The model is described as *full rank* (FR) if $r = \text{rank}([\mathbf{X}' \ \mathbf{R}_x'])' = q$ and otherwise as *less than full rank* (LTFR) if $r < q$.

Optionally, \mathbf{Y}_i may be assumed to follow a jointly Gaussian distribution.

In the present section we assume Σ is unstructured. In later sections we will consider the case in which the $p(p+1)/2$ unique elements of Σ are functions of a smaller number of parameters τ ($k \times 1$), with $k < p(p+1)/2$.

In the following, $\mathbf{t}_j = \text{col}_j(\mathbf{T})$. The bilinear form, $E[\mathbf{y}_{ij} | \text{row}(\mathbf{X}_i), \mathbf{t}_j] = \boldsymbol{\mu}_{ij} = \mathbf{X}_i \mathbf{B} \mathbf{t}_j$, has two interpretations. The first interpretation is that the mean is a linear function of \mathbf{t}_j with regression coefficients ($\boldsymbol{\theta}_i$) that are themselves functions of the characteristics of the participants (e.g., gender), $\boldsymbol{\mu}_{ij} = \boldsymbol{\theta}_i \mathbf{t}_j$ with $\boldsymbol{\theta}_i = \mathbf{X}_i \mathbf{B}$. The second interpretation is that the mean is a linear function of \mathbf{X}_i with regression coefficients ($\boldsymbol{\psi}_j$) which are themselves functions of dose or time, etc., $\boldsymbol{\mu}_{ij} = \mathbf{X}_i \boldsymbol{\psi}_j$ with $\boldsymbol{\psi}_j = \mathbf{B} \mathbf{t}_j$. The concept of statistical interaction unifies the two interpretations. In particular, the magnitudes of the main effects and slopes for the variables represented in \mathbf{X} (or \mathbf{T}) depend on values in \mathbf{T} (or \mathbf{X}). The next lemma allows concluding that the GCM represents mean response as a linear function of the qm cross products of the q explanatory variables represented in row \mathbf{X}_i and the m explanatory variables represented in column \mathbf{t}_j . If $\mathbf{X}_i = [x_{i1} \ x_{i2} \ x_{i3}]$ and $\mathbf{t}_j = [1 \ d_j \ d_j^2]'$, then the expected value for participant i on occasion j contains nine cross products, $\{x_{i1}, x_{i1}d_j, x_{i1}d_j^2, x_{i2}, x_{i2}d_j, x_{i2}d_j^2, x_{i3}x_{i3}d_j, x_{i3}d_j^2\}$.

Lemma 4.2 The GCM assumptions have the following implications. For participant i on occasion j the mean response is

$$E(\mathbf{y}_{ij} | \mathbf{X}, \mathbf{T}) = \text{vec}(\mathbf{X}_i \mathbf{B} \mathbf{t}_j) = (\mathbf{X}_i \otimes \mathbf{t}_j) \text{vec}(\mathbf{B}'). \quad (4.6)$$

For participant i the mean response vector for all times together is

$$E(\mathbf{Y}'_i | \mathbf{X}, \mathbf{T}) = \text{vec}[(\mathbf{X}_i \mathbf{B} \mathbf{T})'] = (\mathbf{X}_i \otimes \mathbf{T}') \text{vec}(\mathbf{B}'). \quad (4.7)$$

For all observations simultaneously

$$E[\text{vec}(\mathbf{Y}') | \mathbf{X}, \mathbf{T}] = \text{vec}[(\mathbf{X} \mathbf{B} \mathbf{T})'] = (\mathbf{X} \otimes \mathbf{T}') \text{vec}(\mathbf{B}'), \quad (4.8)$$

which is of the form $\mathbf{X}_* \boldsymbol{\beta}_*$ ($Np \times qm$)($qm \times 1$).

Data analysts often assume the errors follow a Gaussian distribution. As detailed in Chapter 8, writing $\mathbf{Y} \sim \mathcal{N}_{n,m}(\mathbf{M}, \Xi, \Sigma)$ indicates \mathbf{Y} follows a direct-product matrix Gaussian distribution. By definition, Ξ and Σ are symmetric and positive definite or positive semidefinite, and $\text{vec}(\mathbf{Y}') \sim \mathcal{N}_{n,m}[\text{vec}(\mathbf{M}'), \Xi \otimes \Sigma]$.

Definition 4.6 Writing $GCM_{N,p,q,m}(Y_i; X_i BT, \Sigma)$ with Gaussian errors indicates $Y_i' \sim \mathcal{N}_p\{\text{row}_i(\mathbf{X})\mathbf{B}'\}, \Sigma\}$. Equivalently, $Y \sim \mathcal{N}_{N,p}(\mathbf{X}\mathbf{B}, \mathbf{I}_N, \Sigma)$.

Example 4.2 Ordered responses of interest are repeated measures of height in centimeters recorded at ages 2, 3, 4, and 6 years for n boys and n girls. It may be plausible to suppose mean height increases linearly with age in years 2–6. In the context of $E(Y|\mathbf{X}, \mathbf{T}) = \mathbf{X}\mathbf{B}\mathbf{T}$ ($N \times q$)($q \times m$)($m \times p$) there are $n = N/2$ participants per group, with $q = 2$ groups. The $p = 4$ measurements per participant are indexed by age, $\mathbf{d}' = [2\ 3\ 4\ 6]$, and

$$\begin{aligned} \mathbf{X}\mathbf{B}\mathbf{T} &= \begin{bmatrix} \mathbf{1}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_n \end{bmatrix} \begin{bmatrix} \beta_{11} & \beta_{12} \\ \beta_{21} & \beta_{22} \end{bmatrix} \begin{bmatrix} \mathbf{1}'_p \\ \mathbf{d}' \end{bmatrix} \\ &= (N \times q)(q \times m)(m \times p) \end{aligned} \tag{4.9}$$

Here $\text{row}_i(\mathbf{X}) = \mathbf{X}_i = [x_{i1} \ x_{i2}]$ is $[1\ 0]$ for girls and $[0\ 1]$ for boys. Parameters β_{11} and β_{12} are the intercept and slope, respectively, for girls and parameters β_{21} and β_{22} are the intercept and slope, respectively, for boys. In terms of $\mathbf{X}_i = \text{row}_i(\mathbf{X})$ and $\mathbf{t}_j = \text{col}_j(\mathbf{T})$, the mean for participant i on occasion j is $\mu_{ij} = \mathbf{X}_i \mathbf{B} \mathbf{t}_j$. By the lemma,

$$\mu_{ij} = [x_{i1} \ x_{i1}d_j \ x_{i2} \ x_{i2}d_j] \text{vec}(\mathbf{B}'), \tag{4.10}$$

and the mean of the vertical concatenation of all the rows of \mathbf{Y} is

$$\begin{aligned} E[\text{vec}(\mathbf{Y}')|\mathbf{X}, \mathbf{T}] &= \left(\begin{bmatrix} \mathbf{1}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_n \end{bmatrix} \otimes [\mathbf{1}_p \ \mathbf{d}] \right) \text{vec}(\mathbf{B}') \\ &= \begin{bmatrix} \mathbf{1}_{n \cdot p} & \mathbf{d}_* & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1}_{n \cdot p} & \mathbf{d}_* \end{bmatrix} \begin{bmatrix} \beta_{11} \\ \beta_{12} \\ \beta_{21} \\ \beta_{22} \end{bmatrix}, \end{aligned} \tag{4.11}$$

in which $\mathbf{d}_* = (\mathbf{1}_n \otimes \mathbf{d})$ is a column vector. The difference between slopes, $(\beta_{12} - \beta_{22})$, is an example of gender-by-age interaction. The mean for girls, $\mu_{ij} = \beta_{11} + \beta_{12}d_j$, and the mean for boys, $\mu_{ij} = \beta_{21} + \beta_{22}d_j$, are necessarily of the same functional form because the within-subject design matrix is assumed to be common to all subjects.

4.6 THE RELATIONSHIP OF THE GCM TO THE MULTIVARIATE MODEL

The $GCM_{N,p,s,q}(Y_i T; X_i B T, T' \Sigma T)$ can be interpreted as a transformed multivariate linear model. The corresponding model for all of the data is

$$\mathbf{Y}\mathbf{T} = \mathbf{X}\mathbf{B}\mathbf{T} + \mathbf{E}\mathbf{T}. \tag{4.12}$$

For a single independent sampling unit

$$Y_i T = X_i B T + E_i T \quad (4.13)$$

$$Y_{iT} = X_i B_{iT} + E_{iT}. \quad (4.14)$$

To avoid over parameterization the design matrix T ($m \times p$) must have $m \leq p$. Both sides of the equation

$$E(Y|X, T) = \begin{matrix} & & XBT \\ N \times p & (N \times q)(q \times m)(m \times p) & \end{matrix} \quad (4.15)$$

can be postmultiplied by a generalized inverse such as $T^+ = T'(TT')^{-1}$ ($p \times m$) or $T^- = V^{-1}T'(TV^{-1}T')^{-1}$, in which V^{-1} is arbitrary, $p \times p$, and nonsingular. Doing so gives

$$E[Y_i T'(TT')^{-1} | X, T] = X_i B \quad (4.16)$$

and

$$E[Y_i V^{-1} T'(TV^{-1}T')^{-1} | X, T] = X_i B. \quad (4.17)$$

The right-hand side, $\theta_i = X_i B$, is invariant to the choice of generalized inverse. The vector on the left, $\hat{\theta}'_i = (TV^{-1}T')^{-1}TV^{-1}Y'_i$, is easily recognized as being a weighted least squares estimator (or unweighted if $V = I$) for the subject-specific model $E(Y'_i | T) = \theta'_i T$. If $m = p$, then $T^- = T^{-1}$ ($p \times p$) and

$$E(Y_i T^{-1} | X, T) = \begin{matrix} & X_i B \\ 1 \times p & (1 \times q)(q \times p) \end{matrix}. \quad (4.18)$$

In the simple example of heights measured among boys and girls, measurements were made on four occasions and growth was assumed to be a linear function of age. It is possible to include additional terms in the within-subject regression equation, such as a quadratic term. However, it was assumed that higher order terms are not needed. The assumption constrains the boys' mean, $\mu_{ij} = \beta_{11} + \beta_{12}d_j + \beta_{13}d_j^2 + \beta_{14}d_j^3$, with $\beta_{13} = 0$ and $\beta_{14} = 0$. The same constraint applies for the girls. The general notation for the constrained representation is

$$\begin{aligned} E(Y|X, T) &= XBT \\ &= X [B_1 \ B_2] \begin{bmatrix} T_1 \\ T_2 \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \beta_{11} & \beta_{12} & \beta_{13} & \beta_{14} \\ \beta_{21} & \beta_{22} & \beta_{23} & \beta_{24} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 3 & 4 & 6 \\ 2^2 & 3^2 & 4^2 & 6^2 \\ 2^3 & 3^3 & 4^3 & 6^3 \end{bmatrix}, \quad (4.19) \end{aligned}$$

constrained by $B_2 = \mathbf{0}$ with dimensions B_1 ($q \times m$) and B_2 $q \times (p - m)$. Here $m = 2$ and T is square and full rank. Also

$$\begin{aligned}
 E(\mathbf{Y}\mathbf{T}^{-1}|\mathbf{X},\mathbf{T}) &= \mathbf{X}\mathbf{B} \\
 &= \mathbf{X}[\mathbf{B}_1 \ \mathbf{B}_2],
 \end{aligned}
 \tag{4.20}$$

constrained by $\mathbf{B}_2 = \mathbf{0}$. The form suggests partitioning the left-hand side as $\mathbf{Y}\mathbf{T}^{-1} = [\mathbf{Y}_{*1} \ \mathbf{Y}_{*2}]$, with $E(\mathbf{Y}_{*1}|\mathbf{X},\mathbf{T}) = \mathbf{X}\mathbf{B}_1$ while $E(\mathbf{Y}_{*2}|\mathbf{X},\mathbf{T}) = \mathbf{0}$. Rao (1965) first proposed the reduction of the GCM to an ordinary multivariate GLM.

The GCM is a multivariate GLM constrained by $p - m$ linear restrictions on the regression parameter matrix \mathbf{B} . The unconstrained model is denoted $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}_{\mathbf{Y}|\mathbf{X}}, \Sigma_{\mathbf{Y}|\mathbf{X}})$ with $\mathbf{B}_{\mathbf{Y}|\mathbf{X}} (q \times p)$. The assumption of a full rank within-subject design matrix $\mathbf{T} (q \times m)$ with $m \leq p$ defines the linear constraints $\mathbf{B}_{\mathbf{Y}|\mathbf{X}} = \mathbf{B}\mathbf{T}$. If $m = p$, then \mathbf{T} is a square nonsingular matrix and the number of linear restrictions is zero.

Example 4.3 An example is given by

$$\begin{aligned}
 \mathbf{T} &= \{t_{ij} : t_{ij} = d_j^{i-1}\} \\
 &= \begin{bmatrix} 1 & 1 & 1 & 1 \\ d_1 & d_2 & d_3 & d_4 \\ d_1^2 & d_2^2 & d_3^2 & d_4^2 \\ d_1^3 & d_2^3 & d_3^3 & d_4^3 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 3 & 4 & 6 \\ 2^2 & 3^2 & 4^2 & 6^2 \\ 2^3 & 3^3 & 4^3 & 6^3 \end{bmatrix}.
 \end{aligned}
 \tag{4.21}$$

In turn,

$$E(\mathbf{Y}|\mathbf{X},\mathbf{T}) = \begin{bmatrix} \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} \end{bmatrix} \begin{bmatrix} \beta_{11} & \beta_{12} & \beta_{13} & \beta_{14} \\ \beta_{21} & \beta_{22} & \beta_{23} & \beta_{24} \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 \\ 2 & 3 & 4 & 6 \\ 2^2 & 3^2 & 4^2 & 6^2 \\ 2^3 & 3^3 & 4^3 & 6^3 \end{bmatrix}.
 \tag{4.22}$$

Usually some columns of \mathbf{B} will be assumed to be zero. If \mathbf{T} represents polynomial regression of y_{ij} on d_j , the last few columns of \mathbf{B} correspond to the highest order polynomial terms. The constraint on \mathbf{B} induces a partitioning of both \mathbf{B} and \mathbf{T} , with $\mathbf{B}_2 = \mathbf{0} [(q \times (p - m))]$ and $\mathbf{B}_1 (q \times m)$ is the set of nonzero columns. The constraint $\mathbf{B}_2 = \mathbf{0}$ corresponds to omitting some of the rows of \mathbf{T} from the model, which is conceptually no different from deciding to omit some columns of \mathbf{X} . Usually \mathbf{B}_1 is the first few columns of \mathbf{B} , but any columns of \mathbf{B} may be required to be zero. Since the rows of \mathbf{T} and the columns of \mathbf{B} can be permuted, there is no loss in generality in using the notation given above. When some columns of \mathbf{B} are required to be zero, then only the q nonzero columns, \mathbf{B}_1 , must be estimated ($2 \leq m \leq p$). The model can be represented by $\text{GCM}_{N,p,q,m}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}_1\mathbf{T}_1, \Sigma)$ with $\mathbf{B}_1 (q \times m)$ being a subset of the columns of $\mathbf{B} (q \times p)$.

Definition 4.7 In the GCM, the mean response as a function of dose d or time is of the form $\mu(d; \mathbf{C}) = \mathbf{C}\mathbf{B}\mathbf{u}_d$ ($1 \times q \times m \times 1$) in which

$\mathbf{u}_d = [1 \ d \ \cdots \ d^{m-1}]'$ ($m \times 1$). In the special case of $d = d_j$ and $\mathbf{C} = \mathbf{X}_i$ we have $\mu(d_j; \mathbf{X}_i) = E(y_{ij} | \mathbf{X}_i, \mathbf{t}_j) = \mathbf{X}_i \mathbf{B} \mathbf{t}_j$. A set of coordinates $\{[d, \mu(d; \mathbf{C})] : d \in [d_1, d_p]\}$ defines a *growth curve*.

4.7 MIXED, HIERARCHICAL, AND RELATED MODELS

The general linear mixed model encompasses the most general range of models considered in detail in the present book. In terms of expressions for population parameters, many other models mentioned can be thought of as special cases. The special relationships have led to the widespread misconception that similar special case relationships hold for estimates and tests. However, appropriate tests occur as special cases for only a very limited range of models. The mixed model has mostly approximate results and very few exact results for inference. The exact results for estimates and tests for a (univariate) $GLM_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ with Gaussian errors do occur automatically as special cases of mixed model approximations. In contrast, for the $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors, only the maximum likelihood estimates occur automatically as special cases of mixed model results. None of the commonly used tests in mixed models correspond to multivariate model tests, except asymptotically or in the special case of a univariate model.

The approximate theory of mixed models also applies to many models described as “hierarchical,” which allow additional freedom in specifying the random components (in contrast to general linear mixed models). Raudenbush and Bryk (2002) provided a book-length treatment. Mixed model theory also appears to encompass and apply to a wide range of “state-space” models (Billio and Monfort, 1998, provided an example).

CHAPTER 5

The Linear Mixed Model

5.1 MOTIVATION

As discussed in earlier chapters, the general linear mixed model allows missing or mistimed data as well as repeated covariates. The approach also allows specifying covariance structures as a function of a small number of parameters. In most uses, the model implicitly assumes commensurate data (all measured in the same units). Most, but certainly not all, applications involve repeated measures.

Many different classes of models have been described as “mixed models.” The term dates back to early developments in ANOVA. The simplest ANOVA design involves one factor, a categorical predictor with G levels defining G groups, and $n_g = N/G$ independent sampling units in each group (cell). Classical less-than-full-rank coding led to writing the model as a scalar equation, with $i \in \{1, 2, \dots, n_g\}$ and $g \in \{1, 2, \dots, G\}$:

$$y_{ig} = \mu + \alpha_g + e_{ig}. \quad (5.1)$$

Here μ and $\{\alpha_g\}$ are fixed and unknown finite constants (parameters) which characterize the means and $e_{ig} \sim \mathcal{N}(0, \sigma_e^2)$, with σ_e^2 a fixed and unknown finite constant, the variance. The nature of $\{\alpha_g\}$ led to describing the predictor variable as a “fixed effect” and the model as a fixed effect model. In contrast, a “random effect” model assumes $\{\alpha_g\}$ are randomly selected from an infinite population, with independent and identically distributed $\alpha_g \sim \mathcal{N}(0, \sigma_g^2)$ independent of $\{e_{ig}\}$. Here μ alone represents the mean, while $\alpha_g + e_{ig}$ represents total variance in terms of two components. More generally, a model with two or more fixed effects was described as a fixed effects model, while a model with two or more random effects was described as a random effects model. In turn, a model with one or more fixed effects and one or more random effects was referred to as a mixed effects model.

The terms random effect, fixed effect, and mixed effect do not always clearly convey the underlying simple structure of a mixed model to readers not intimately familiar with the theory. We prefer to discuss the parameters of a mixed model in terms of two separate components: the model for the means and the model for the covariance. With Gaussian distributions, specifying the first two moments fully determines the distributions and therefore implicitly all derived properties. The

approach allows couching most discussions about the mixed model in terms of the simple concepts of means, variances and correlations.

Mixed models are a useful tool for longitudinal data exhibiting missing values, inconsistently timed observations, or mistimed observations. Incomplete data put the analyst at a disadvantage because critically important information is missing. Analysis cannot proceed unless the missing information is replaced by assumptions. The mixed model approach presents intuitively appealing assumptions and has some procedures for inference that work well at least with large and some moderate sample sizes.

The approach involves building a model for the expected values of the data and also building a model for the covariances of the data. The linear model for the mean allows the flexibility of using polynomials, trigonometric functions, and regression splines, among others. Both linear and nonlinear models for the covariance structure are useful.

We begin by considering an extremely general class of linear models. As stated, the model allows specifying a particular covariance model in a variety of ways. The generality allows ambiguity without adding further constraints. However, it serves well as the basis for special cases of interest.

The interested reader seeking additional details may wish to consult any of a number of excellent book-length treatments centered on mixed models, including Vonesh and Chinchilli (1997), Khuri, Mathew, and Sinha (1998), Verbeke and Molenberghs (2000), and Demidenko (2004). Timm (2002) discussed the mixed model as an alternative to a wide variety of multivariate methods.

5.2 DEFINITION OF THE MIXED MODEL

Definition 5.1 A *general linear mixed model* will be indicated by $LMM_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta}, \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}'_i + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)]$ and includes the following assumptions. When no clarity will be lost, the model may be abbreviated $LMM_{N,p,q,m}(\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta}, \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}\mathbf{Z}'_i + \boldsymbol{\Sigma}_{e_i})$.

1. For $i \in \{1, 2, \dots, N\}$,
 - (a) the N random vectors, $\{\mathbf{e}_i\}$, are $p_i \times 1$ and mutually independent,
 - (b) the N random vectors, $\{\mathbf{d}_i\}$ are $m \times 1$ and mutually independent, and
 - (c) the $\{\mathbf{e}_i\}$ and $\{\mathbf{d}_i\}$ are all mutually independent.
2. Each \mathbf{X}_i , the $p_i \times q$ *expected value design matrix* for independent sampling unit i , is fixed and known without appreciable error, conditional on knowing the sampling units, for data analysis.
3. Elements of $\boldsymbol{\beta}$ ($q \times 1$) are fixed and unknown and often regression coefficients or means.
4. Each \mathbf{Z}_i , the $p_i \times m$ *covariance design matrix* for independent sampling unit i , is fixed and known without appreciable error, conditional on knowing the sampling units, for data analysis.

5. For $i \in \{1, 2, \dots, N\}$, (a) $E(\mathbf{e}_i) = \mathbf{0}$ and (b) $E(\mathbf{d}_i) = \mathbf{0}$.
 6. With $\boldsymbol{\tau}' = [\boldsymbol{\tau}'_d \boldsymbol{\tau}'_e]$ fixed and unknown, \mathbf{d}_i and \mathbf{e}_i have a finite, fixed, and unknown (joint) covariance matrix which is either positive definite or positive semidefinite:

$$\mathcal{V}\left(\begin{bmatrix} \mathbf{d}_i \\ \mathbf{e}_i \end{bmatrix}\right) = \begin{bmatrix} \boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d) & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e) \end{bmatrix}. \quad (5.2)$$

Elements of $\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)$ are twice differentiable functions of $\boldsymbol{\tau}_d$, a vector of no more than $m(m+1)/2$ covariance parameters. Elements of $\boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)$ are twice differentiable functions of $\boldsymbol{\tau}_e$, a vector of no more than $\max_i [p_i(p_i+1)/2]$ covariance parameters.

7. The $p_i \times 1$ response vector \mathbf{y}_i is expressed as $\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i$ with

$$E(\mathbf{y}_i) = \mathbf{X}_i\boldsymbol{\beta} \quad (5.3)$$

and fixed, unknown, and positive definite or positive semidefinite covariance

$$\mathcal{V}(\mathbf{y}_i) = \boldsymbol{\Sigma}_i(\boldsymbol{\tau}) = \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}'_i + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e). \quad (5.4)$$

When no clarity will be lost, the covariance model may be abbreviated $\boldsymbol{\Sigma}_i$.

Writing $\text{LMM}_{N,p,q,m}(\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}_x\boldsymbol{\beta} = \mathbf{a}, \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}\mathbf{Z}'_i + \boldsymbol{\Sigma}_{e_i} | \mathbf{R}_d\mathbf{d}_i = \mathbf{0})$ specifies explicit restrictions on parameters in $\boldsymbol{\beta}$ through the fixed and known constants \mathbf{R}_x and \mathbf{a} and explicit restrictions on \mathbf{d}_i through the fixed and known constant \mathbf{R}_d .

The definition contains what may be described as the “approximate generalized least squares” assumptions. In combination with mild restrictions on dimensions and ranks of $\{\mathbf{X}_i, \boldsymbol{\Sigma}_i\}$, they guarantee the existence of estimators for $\boldsymbol{\beta}$ and $\boldsymbol{\tau}$ which satisfy an approximate generalized least squares criterion. It is very important to recognize that no particular distribution has been specified for any random variable. Only the rather modest requirement of finite second (and implicitly first) moments is made.

The model definition specifies three components: the response vector for the independent sampling unit, the mean of the response vector, and the covariance of the response vector. In turn, the covariance of the response consists of two components, corresponding to the two unobservable random vectors \mathbf{d}_i and \mathbf{e}_i .

The data may be “stacked” to represent a combined model. With $n = \sum_{i=1}^N p_i$, it is often convenient to write $\mathbf{y}'_s = [\mathbf{y}'_1 \mathbf{y}'_2 \cdots \mathbf{y}'_N]$, which implies \mathbf{y}_s is $n \times 1$. Similarly, $\mathbf{d}'_s = [\mathbf{d}'_1 \mathbf{d}'_2 \cdots \mathbf{d}'_N]$ implies \mathbf{d}_s is $Nm \times 1$ and $\mathbf{e}'_s = [\mathbf{e}'_1 \mathbf{e}'_2 \cdots \mathbf{e}'_N]$ implies \mathbf{e}_s is $n \times 1$. In turn, concatenation of the predictor matrices gives

$$\mathbf{X}_s = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}, \quad (5.5)$$

an $n \times q$ matrix. In contrast,

$$\mathbf{Z}_s = \begin{bmatrix} \mathbf{Z}_1 & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{Z}_2 & & \vdots \\ \vdots & & \ddots & \mathbf{0} \\ \mathbf{0} & \cdots & \mathbf{0} & \mathbf{Z}_N \end{bmatrix} \tag{5.6}$$

is an $n \times (Nm)$ block diagonal matrix. Necessarily the rows of \mathbf{y}_s , \mathbf{d}_s , \mathbf{e}_s , \mathbf{X}_s , and \mathbf{Z}_s must be sorted in the same order (both within and between independent sampling units) for the following equation, the stacked data model, to be valid:

$$\mathbf{y}_s = \mathbf{X}_s\boldsymbol{\beta} + \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s. \tag{5.7}$$

Consequently

$$E(\mathbf{y}_s) = \mathbf{X}_s\boldsymbol{\beta} \tag{5.8}$$

and, with $\Sigma_{ds} = \bigoplus_{i=1}^N \Sigma_{di}$,

$$\begin{aligned} \mathcal{V}(\mathbf{y}_s) &= \Sigma_s = \bigoplus_{i=1}^N (\mathbf{Z}_i \Sigma_{di} \mathbf{Z}'_i + \Sigma_{ei}) \\ &= \bigoplus_{i=1}^N (\mathbf{Z}_i \Sigma_{di} \mathbf{Z}'_i) + \bigoplus_{i=1}^N \Sigma_{ei} \\ &= \mathbf{Z}_s \Sigma_{ds} \mathbf{Z}'_s + \Sigma_{es}. \end{aligned} \tag{5.9}$$

The model $\mathbf{y}_s = \mathbf{X}_s\boldsymbol{\beta} + \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s$ expresses the observations as a function of three terms. The term $\mathbf{X}_s\boldsymbol{\beta}$ describes the fixed contribution of the population, conditional on the predictor values (which often define subpopulations). Each row of $\mathbf{Z}_s\mathbf{d}_s$ describes a random deviation from the population value due to observing a particular person (ISU). Each row of \mathbf{e}_s describes an additional and distinct random deviation due to observing a particular person on a particular occasion. In summary, a mixed model expresses an observation as a *subpopulation mean* plus a *random deviation due to person* plus an additional *random deviation due to occasion*.

Alternately, the response vector may be expressed in terms of one purely fixed and one purely random vector. If $\mathbf{e}_{+s} = \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s$, then

$$\begin{aligned} \mathbf{y}_s &= \mathbf{X}_s\boldsymbol{\beta} + \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s \\ &= \mathbf{X}_s\boldsymbol{\beta} + \mathbf{e}_{+s} \\ &= \text{fixed} + \text{random} \\ &\quad \uparrow \qquad \qquad \uparrow \\ &\quad E(\mathbf{y}_s) \qquad \mathcal{V}(\mathbf{y}_s) \\ &\quad \text{model} \qquad \text{model} \end{aligned} \tag{5.10}$$

The fixed component $\mathbf{X}_s\boldsymbol{\beta}$ completely determines the mean, the first moment, of \mathbf{y}_s and has no effect on any variance or covariance. The random component \mathbf{e}_{+s} completely determines all variances and covariances, the second moments, and has no effect on the mean. The model for a particular independent sampling unit

naturally follows the same pattern:

$$\begin{aligned}
 \mathbf{y}_i &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i \\
 &= \mathbf{X}_i\boldsymbol{\beta} + \mathbf{e}_{+i} \\
 &= \text{fixed} + \text{random}.
 \end{aligned}$$

$$\begin{array}{ccc}
 & \uparrow & \uparrow \\
 & \text{E}(\mathbf{y}_i) & \mathcal{V}(\mathbf{y}_i) \\
 & \text{model} & \text{model}
 \end{array}
 \tag{5.11}$$

As in univariate and multivariate linear models, a less-than-full-rank design matrix \mathbf{X}_s prevents unique identification and unbiased estimation of $\boldsymbol{\beta}$. Side conditions must be imposed to avoid the problem. We shall usually refer to the side conditions as “restrictions” or “constraints.” As the model is defined in Definition 1.1, side conditions also must be imposed on the components of Σ_i to uniquely define the covariance model parameters and estimates. Given appropriate side conditions, modern computing tools often make it straightforward to find estimates for $\boldsymbol{\beta}$ and $\boldsymbol{\tau}$ which satisfy an approximate least squares criterion or an iterated approximate least squares criterion. Except for special cases or in large samples, the estimators have few optimal properties. Kackar and Harville (1984) proved that $\hat{\boldsymbol{\beta}}$ from iterated approximate least squares (sometimes called estimated generalized least squares, among other names) is unbiased for a Gaussian (or any other symmetric) distribution. However, covariance parameter estimators, at least in small samples, typically have substantial bias. Littell (2003) provided an excellent overview.

With additional restrictions on the covariance parameters $\boldsymbol{\tau} = [\boldsymbol{\tau}'_d \boldsymbol{\tau}'_e]'$ the model defined in Definition 1.1 becomes sufficiently well-defined to allow computing parameter estimators with reasonable properties. Typically data analysts greatly simplify or completely eliminate either $\boldsymbol{\tau}_d$ or $\boldsymbol{\tau}_e$. As an example, assuming response vector i has a first-order autocorrelation covariance pattern requires the j, k element to be $\{\Sigma_i(\boldsymbol{\tau})\}_{j,k} = \sigma^2\rho^{|t_j-t_k|}$. Choosing $\mathbf{Z}_i = \mathbf{0}$, $\Sigma_{ei}(\boldsymbol{\tau}_e) = \{\sigma^2\rho^{|t_j-t_k|}\}$, and $\boldsymbol{\tau}_e = [\sigma^2 \rho]'$ achieves the desired pattern. Doing so expresses the variances and covariances as a nonlinear function of the two parameters σ^2 and ρ . Implicitly, the observations for sampling unit i follow a stationary time series. For nonstationary processes an inherently linear model for the covariances is often assumed. In particular, specifying $\{\mathbf{G}_{ik}\}$ as known constant matrices allows writing

$$\Sigma_i(\boldsymbol{\tau}) = \sum_{k=1}^t \tau_k (\mathbf{Z}_i \mathbf{G}_{ik} \mathbf{Z}'_i) + \sigma^2 \mathbf{I}_{p_i}.
 \tag{5.12}$$

With such a structure $\Sigma_i(\boldsymbol{\tau})$ is a linear function of $t + 1$ unknown parameters.

Example 5.1 The ambiguity in representing the covariance model may be illustrated quite easily. The response vector (and the purely random part of the mixed model) will have a compound symmetric covariance structure if $\mathcal{V}(\mathbf{e}_i) = \sigma^2 [\mathbf{1}_{p_i} \mathbf{1}'_{p_i} \rho + \mathbf{I}_{p_i} (1-\rho)]$ and $\mathbf{Z}_i = \mathbf{0}$ [$\mathbf{Z}_i = \mathbf{0}$ has exactly the same effect as

assuming $\mathcal{V}(\mathbf{d}_i) = \mathbf{0}$]. In either case $\mathcal{V}(\mathbf{y}_i) = \mathcal{V}(\mathbf{e}_i)$. Alternately, choosing $\mathbf{Z}_i = \mathbf{1}_{p_i}$, $\mathcal{V}(\mathbf{d}_i) = \sigma^2 \rho$, and $\mathcal{V}(\mathbf{e}_i) = \sigma^2(1-\rho)\mathbf{I}_{p_i}$ gives exactly the same $\mathcal{V}(\mathbf{y}_i)$, but with a different $\mathcal{V}(\mathbf{e}_i)$.

Most presentations of the general linear mixed model include the assumption that the random variables follow a Gaussian distribution. However, in contrast to the univariate and multivariate linear models, the assumption does not lead to closed formed expressions for estimates.

Definition 5.2 Writing

LMM $_{N,p_i,q,m}(\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}_x\boldsymbol{\beta} = \mathbf{a}, \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}_i' + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e) | \mathbf{R}_d\mathbf{d}_i = \mathbf{0})$
with Gaussian errors indicates

$$\begin{bmatrix} \mathbf{d}_i \\ \mathbf{e}_i \end{bmatrix} \sim \mathcal{N}_{m+p_i} \left\{ \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d) & \mathbf{0} \\ \mathbf{0} & \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e) \end{bmatrix} \right\}, \tag{5.13}$$

which is equivalent to

$$\begin{bmatrix} \mathbf{d}_s \\ \mathbf{e}_s \end{bmatrix} \sim \mathcal{N}_{2n} \left\{ \begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} \bigoplus_{i=1}^N \boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d) & \mathbf{0} \\ \mathbf{0} & \bigoplus_{i=1}^N \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e) \end{bmatrix} \right\}. \tag{5.14}$$

Recalling $\boldsymbol{\Sigma}_i(\boldsymbol{\tau}) = \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}_i' + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)$, it follows that

$$\mathbf{y}_i \sim \mathcal{N}_{p_i}[\mathbf{X}_i\boldsymbol{\beta}, \boldsymbol{\Sigma}_i(\boldsymbol{\tau})]. \tag{5.15}$$

Equivalently,

$$\mathbf{y}_s \sim \mathcal{N}_n \left[\mathbf{X}_s\boldsymbol{\beta}, \bigoplus_{i=1}^N \boldsymbol{\Sigma}_i(\boldsymbol{\tau}) \right]. \tag{5.16}$$

Some mixed models, members of a class often referred to as components of variance models, assume compound symmetric covariance among *all* observations, which implies having only one independent sampling unit. The special properties of compound symmetry allow using exact weighted least squares methods to transform the model to one with completely independent observations (and some heterogeneity).

The mnemonic HILE Gauss must be interpreted carefully in the mixed model setting because allowing p_i to vary causes $\boldsymbol{\Sigma}_i(\boldsymbol{\tau})$ to vary. *Independence* of sampling units remains the cornerstone (even given the special handling required for some components of variance models). In turn, *linearity* of the response expected value (mean) as a function of the parameters also holds, as does the assumption of finite second moments (*existence*). However, describing the model

as having *homogeneity* of second moments must be interpreted to indicate that a single covariance *model* holds for all sampling units, while $\Sigma_i(\boldsymbol{\tau}) = \Sigma_j(\boldsymbol{\tau})$ only if $\mathbf{Z}_i = \mathbf{Z}_j$ and $p_i = p_j$.

5.3 DISTRIBUTION-FREE AND NONITERATIVE ESTIMATES

As always in the history of mixed models, computing difficulties often dominate theory and practice. Although current computing hardware speed and associated software have virtually eliminated some problems, we still face many issues considered by Henderson (1953). Method-of-moments and MINQUE (minimum variance quadratic unbiased estimation) provide noniterative estimates (Searle, Casella, and McCulloch, 1992). Although the methods have many good properties, they are currently much less popular than likelihood methods.

5.4 GAUSSIAN LIKELIHOOD AND ITERATIVE ESTIMATES

The nearly ubiquitous use of the Gaussian assumption leads most data analysts to seek either maximum likelihood (ML) estimates or *restricted maximum likelihood* (REML) estimates (Chapter 14 has some details). The joint log likelihood is, with $n = \sum_{i=1}^N p_i$, $\mathbf{e}'_{+i} = \mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}$ and $\Sigma_i(\boldsymbol{\tau})$ abbreviated as Σ_i ,

$$-2\log L(\boldsymbol{\beta}, \boldsymbol{\tau}) = n\log(2\pi) - \frac{1}{2} \sum_{i=1}^N [\log|\Sigma_i| + (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})' \Sigma_i^{-1} (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})]. \quad (5.17)$$

Iterative methods must be used to solve the system of equations, which are nonlinear in the parameters $\{\boldsymbol{\beta}, \boldsymbol{\tau}\}$. Without the Gaussian assumption, the resulting estimates satisfy the iterated approximate weighted least squares criterion (ITAWLS). Either with or without the Gaussian assumption, the resulting estimates are biased in small samples. As mentioned earlier, the estimate of $\boldsymbol{\beta}$ is unbiased, while the estimate of $\boldsymbol{\tau}$ (and $\{\Sigma_i\}$) typically has substantial bias in small samples.

REML estimates of $\boldsymbol{\tau}$ (and $\{\Sigma_i\}$) have less bias than ML estimates. REML estimates arise from maximizing the reduced profile log-likelihood equation, based on $\hat{\mathbf{e}}'_{+i} = \mathbf{y}_i - \mathbf{X}_i\hat{\boldsymbol{\beta}}$:

$$-2\log \hat{L}_{\text{REML}}(\boldsymbol{\tau}) = (n-q)\log(2\pi) + \sum_{i=1}^N (\log|\Sigma_i| + \hat{\mathbf{e}}'_{+i} \Sigma_i^{-1} \hat{\mathbf{e}}_{+i}) + \log \left| \sum_{i=1}^N \mathbf{X}'_i \Sigma_i^{-1} \mathbf{X}_i \right|. \quad (5.18)$$

Although we leave the details to Chapter 14, it is worth observing the following. In the univariate linear model, a special case of the mixed model, the maximum likelihood estimate of the error variance, $\hat{\sigma}^2 = \mathbf{y}'[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}/N$, has

an expected value smaller than σ^2 . In contrast, $\hat{\sigma}^2 = \tilde{\sigma}^2 N / [N - \text{rank}(\mathbf{X})]$ is the unbiased REML estimate of σ^2 .

Estimates of both the primary parameters and a variety of secondary parameters are usually desired. Interest may center on either estimates of population location parameters (elements of $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$) which apply to all sampling units or estimates of properties particular to a specific independent sampling unit. We leave further discussion to the more detailed treatment of estimation and testing in later chapters.

5.5 TESTS ABOUT $\boldsymbol{\beta}$ (MEANS, FIXED EFFECTS)

In the context of the mixed model, data analysts often wish to test a hypothesis of the form

$$\begin{aligned} H_0 : \mathbf{C}\boldsymbol{\beta} &= \boldsymbol{\theta}_0 \\ H_0 : \boldsymbol{\theta} &= \boldsymbol{\theta}_0. \end{aligned} \quad (5.19)$$

Only low order approximate tests have been described. Not surprisingly, the tests often perform very poorly in small samples.

Given the assumption of Gaussian data and the strong parallels of the forms to the univariate linear model, it is straightforward to define the form of the likelihood ratio test statistic. Computing the statistic requires iterative calculations to successfully fit two models, the full model and the constrained model, at least one of which must be false. In practice, the false model is less likely to converge. If both models converge, then the test is well defined. The difficulty with the likelihood ratio test lies in finding an adequate approximation to the distribution of the test statistic. The approximation $-2\log[L(\hat{\boldsymbol{\theta}}|\mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0)/L(\hat{\boldsymbol{\theta}})] \sim \chi^2(a)$ with $\boldsymbol{\theta}$ of dimension $a \times 1$ relies on the log likelihood being approximately quadratic in shape. From the perspective of a Taylor series expansion for the test statistic, the approximation is less than first order (it matches the first moment only asymptotically). The inaccuracy arises from ignoring the variability due to replacing $\boldsymbol{\Sigma}_i$ by $\hat{\boldsymbol{\Sigma}}_i$.

It is also possible to create tests based on model comparisons with estimates based on REML estimates. Again, little is known about the resulting distributions.

A variety of alternative tests have been suggested based on assuming the distribution of the test statistic may be approximated by an F distribution (sometimes referred to as a Wald type test). Such tests can be based on ML or REML estimates. A key advantage lies in only needing to fit a single model (which ideally provides an essentially correct model). With estimates $\hat{\boldsymbol{\theta}} = \mathbf{C}\hat{\boldsymbol{\beta}}$ and

$$\hat{\boldsymbol{\Sigma}}_s = \bigoplus_{i=1}^N \hat{\boldsymbol{\Sigma}}_i \quad (5.20)$$

and \mathbf{C} of dimension $a \times q$ and rank a ,

$$\mathcal{V}(\widehat{\boldsymbol{\theta}}) = \mathcal{V}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) = \mathbf{C}(\mathbf{X}'_s \widehat{\boldsymbol{\Sigma}}_s^{-1} \mathbf{X}_s)^{-1} \mathbf{C}' \tag{5.21}$$

and

$$F_m = (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'_s \widehat{\boldsymbol{\Sigma}}_s^{-1} \mathbf{X}_s)^{-1} \mathbf{C}']^{-1} (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) / a. \tag{5.22}$$

Among the widely available methods, the Kenward and Roger (1997) method appears to provide the most accurate test size. However, room for substantial improvement remains (Schaalje, McBride, and Fellingham, 2003). The Kenward Roger method starts with REML estimates and then creates an improved estimate of $\boldsymbol{\Sigma}_s$. Sample values are used to estimate a scale parameter λ and a degrees-of-freedom parameter ν with

$$F_m \widehat{\lambda} \sim F(a, \nu). \tag{5.23}$$

5.6 TESTS OF COVARIANCE PARAMETERS, $\boldsymbol{\tau}$ (RANDOM EFFECTS)

As in the univariate and multivariate linear models, a data analyst may wish to test hypotheses about variance or covariance parameters. The likelihood ratio test provides a reasonable approach. A variety of approximate F approaches also have been proposed, based on extending consideration to

$$\begin{bmatrix} \boldsymbol{\theta}_\beta \\ \boldsymbol{\theta}_\tau \end{bmatrix} = \begin{bmatrix} \mathbf{C}_\beta & \mathbf{C}_\tau \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta} \\ \boldsymbol{\tau} \end{bmatrix}. \tag{5.24}$$

Such tests have not received much attention. Some work of a similar nature has been done in univariate and multivariate linear models.

EXERCISES

5.1 Clearly specify values for every dimension and parameter which reduce a $\text{LMM}_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i \boldsymbol{\beta}, \mathbf{Z}_i \boldsymbol{\Sigma}_{di}(\boldsymbol{\tau}_d) \mathbf{Z}'_i + \boldsymbol{\Sigma}_{ei}(\boldsymbol{\tau}_e)]$ with Gaussian errors to a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ with Gaussian errors. Some choices will not be unique.

5.2 For a $\text{GLM}_{N,p,q}(\mathbf{Y}; \mathbf{X}; \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors, the model may be written as $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$. Assume $\mathbf{u}_s = \text{vec}(\mathbf{Y}')$, and $\mathbf{u}_i = [\text{row}_i(\mathbf{Y})]'$.

5.2.1 Clearly specify values for every dimension and parameter which specify a $\text{LMM}_{N,p,q,m}[\mathbf{u}_j; \mathbf{X}_{ju} \boldsymbol{\beta}_u, \mathbf{Z}_j \boldsymbol{\Sigma}_{dj}(\boldsymbol{\tau}_d) \mathbf{Z}'_j + \boldsymbol{\Sigma}_{ej}(\boldsymbol{\tau}_e)]$ with Gaussian errors. Some choices will not be unique.

5.2.2 Clearly specify (and simplify the expressions when possible) all matrices (including dimensions and pattern of elements) in the stacked-data form $\mathbf{u}_s = \mathbf{X}_s \boldsymbol{\beta}_u + \mathbf{e}_u$ corresponding to the choices you made in 5.2.1.

5.3 Consider the magnetic resonance imaging (MRI) data used for the Chapter 3 exercises. Ignoring all data for parenchyma and fat gives a multivariate model with three columns for the three region-of-interest (ROI) measurements.

5.3.1 Briefly specify an appropriate multivariate model, $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors. Use reference cell coding. Specify all dimensions for the particular data in hand, as well as brief definitions of parameter matrix elements.

5.3.2 Assume $\mathbf{Z}_i \equiv \mathbf{0}$ and an unstructured covariance pattern across time. Clearly specify values for every dimension and parameter which specify a $\text{LMM}_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta}, \mathbf{Z}_i\Sigma_{di}(\boldsymbol{\tau}_d)\mathbf{Z}'_i + \Sigma_{ei}(\boldsymbol{\tau}_e)]$ with Gaussian errors appropriate for the data which gives estimates of the same parameters as in the $\text{GLM}()$ model.

5.3.3 Clearly specify (and simplify the expressions where possible) all matrices (including dimensions and pattern of elements) in the stacked data form (for *all* observations on *all* participants) corresponding to the choices you made in 5.3.2.

5.4 In the MRI analysis, the multivariate model applies and therefore should always be used in preference to the mixed model. However, to improve understanding of linear models, the data may be analyzed with a mixed model. For the following questions, we are only interested in analyzing the ROI tissue type (ignore all data for parenchyma and fat).

5.4.1 Briefly explain *why* a multivariate model is preferred to a mixed model here.

5.4.2 Use a “stacked” version of the data (in file P0105) to fit the mixed model from exercise 5.3 using SAS PROC MIXED. You will need to use the CLASS and REPEATED statements in PROC MIXED and specify a compound symmetric covariance matrix. Provide tests for the Time and Diagnosis main effects as well as the Time x Diagnosis interaction (malignant versus benign) which correspond to the test that a multivariate model would provide.

5.4.3 Starting with the P0104 file, use PROC GLM with the REPEATED statement to compute MULTIREP and UNIREP tests of Time, Diagnosis, and Time by Diagnosis.

5.4.4 Compare and discuss the degrees of freedom, F statistics, and p values from exercises 5.4.2 and 5.4.3.

5.4.5 Implement a new version of the PROC GLM analysis and use reference cell coding in lieu of the CLASS statement for coding benign versus malignant. What changes result?

5.4.6 (*Optional, noncredit*) Use LINMOD to compute MULTIREP and UNIREP tests of the same three hypotheses.

5.5 (*Optional, noncredit*) Considering both the ROI and Parenchyma tissue type, we will now consider a factorial analysis that includes Time (three levels), Diagnosis (two levels), and Tissue (two levels). Again, the multivariate model is appropriate and should be used, but conducting a mixed model analysis will illustrate many issues.

5.5.1 (*Optional, noncredit*) Starting with the P0105.sd2 file, create a file of stacked data which includes all covariates and the ROI and Parenchyma response variables. Using SAS PROC MIXED with the CLASS and REPEATED statements, provide tests for the Time, Diagnosis, and Tissue main effects as well as the three two-way and one three-way interaction.

5.5.2 (*Optional, noncredit*) Compare and discuss the degrees of freedom, F statistics, and p values from exercise 5.5.1, with exercise 3.4 from Chapter 3.

CHAPTER 6

Choosing the Form of a Linear Model for Analysis

6.1 THE IMPORTANCE OF UNDERSTANDING DEPENDENCE

The pattern of dependence among observations can have more effect on the validity and quality of a statistical analysis than any other feature. Therefore choosing an analysis must begin with determining the logical properties of the sampling scheme and thereby characterizing the patterns of independence and correlation among observations.

In the simplest case, as in the general linear univariate model, all response values, all observations, are statistically independent. Relatively simple and well-behaved methods are nearly always available for completely independent observations. In the most complicated case, all observations are correlated with each other in idiosyncratic ways. Few methods with desirable properties can be found for such general problems. The following definitions help characterize fundamental properties of sampling schemes central to the choice of an analysis.

Definition 6.1 (a) *Independent* observations have values which are statistically independent.

(b) An *independent sampling unit* (ISU) provides one or more observations such that observations from one unit are statistically independent from any other distinct unit while observations from the same unit may be correlated.

(c) The *observational unit* distinguishes one correlated observation from another within the ISU.

(d) Observing the same variable in two or more instances across time, space, or other dimension within an ISU creates *repeated measures*.

(e) *Commensurate* observations share the same measurement scale and units.

(f) *Multivariate* outcomes arise from a single ISU and therefore are not independent and need not be commensurate.

(g) *Doubly multivariate* outcomes include repeated measures of two or more noncommensurate variables.

Developers of multivariate statistics have progressed by considering a relatively small number of limited and systematic patterns of dependence among observations. A simple taxonomy useful for describing any type of data and pattern of dependence may be developed by considering three diagnostic questions. We consider each separately.

6.2 HOW MANY VARIABLES PER INDEPENDENT SAMPLING UNIT?

Table 6.1 summarizes some dimensions for describing models. In a clinical trial of a new pharmaceutical with random assignment to treatment, usually the individual participants in the trial are the independent sampling units. Each person may have more than one response measured, such as red blood cell count and blood cholesterol level, which leads to two observations per ISU. For the example, a single value of a blood assay represents the observation unit. As a second example, an educator who randomly assigns all children in a classroom to a particular curriculum must treat scores from children within a classroom as correlated.

Table 6.1 How Many Variables?

Number of Responses	Number of Predictors	Model Description
1	1	Univariate
1	Many	Multivariable
Many	1 or many	Multivariate
Many	Many	Multivariate
Repeated	1 or many	Repeated measures

In the latter setting, classroom becomes the ISU and child the observational unit. Recording performance on the same test for every child leads to repeated measures occurring within the classroom (ISU). Alternately, measuring only one child from each classroom once per month for three months also leads to repeated measures. All repeated measures have commensurate observations (all measured in the same units, on the same scale). However, not all sets of commensurate data are analyzed appropriately with repeated-measures models, which usually imply interest in contrasts corresponding to polynomial trends across a repeated-measure dimension. A simple example arises in measuring drug or toxicant levels in a variety of organs in the body of an animal. Although levels in the brain, muscle, and kidney would all be reported in the same units (the data are commensurate), trends across organs have no scientific appeal.

The multivariate profile of drug or toxicant levels for a variety of organs does pique scientific interest. Statisticians usually describe any setting with two or more response variables as multivariate. From the perspective of mathematical and statistical theory, repeated measures merely represent a special case of multivariate

responses. Modeling and testing differences *between* independent sampling units often allow using simple univariate theory, while modeling and testing differences *within* independent sampling units usually require more complex multivariate theory, as do differences involving both within and between dimensions.

Some research involves collecting two or more distinct response variables on two or more occasions. Many longitudinal studies of human development and clinical trials of treatments for chronic diseases have such data, which may be described as doubly multivariate.

In summary, completely independent observations may be contrasted with many patterns of nonindependent observations. The presence of a nonzero correlation provides the simplest way to identify nonindependent observations, because any correlated variables are necessarily dependent (*not* independent). Although uncommon, examples of uncorrelated and nonindependent variables do exist. Hence a lack of correlation does not guarantee independence. However, in the special case of jointly Gaussian variables, the converse does hold: Uncorrelated and jointly Gaussian random variables are necessarily independent.

The emphasis on recognizing the many forms of nonindependence reflects the crucial importance that the pattern of dependence plays in determining the underlying distribution theory and the choices of parameter estimates and tests. Improper analysis can severely bias results for estimation and inference. Furthermore, in sharp contrast to most other assumptions in linear models, access to large samples will not always overcome the problem.

The term “large” is ambiguous in the presence of multivariate data. If a sample contains p observations on each of N independent sampling units, consideration must be given to increasing N alone (the most common meaning), p alone, N and p in a fixed ratio, or N and p in a varying ratio (with many variations).

Fortunately, the theory for multivariate and repeated-measures data coincides for the multivariate general linear model. Only the particulars of the scientific context and goals can indicate which analysis method to choose for a specific application. The practical use and interpretation of multivariate theory varies greatly across applications, in contrast to the theory itself.

6.3 WHAT TYPES OF VARIABLES PLAY A ROLE?

Definition 6.2 (a) *Nominal* scales only define categories or groups of observations.

(b) *Ordinal* scales provide numeric values sufficient only to rank observations.

(c) *Interval* scales provide numeric values with all differences of the same size being equivalent.

(d) *Ratio* scales give numeric values for which ratios of the same size are equivalent.

(e) *Continuous* data may include any sort of interval- and ratio-scale variables.

The type of measurement scale, especially for random components such as error terms, typically plays a major role in determining the choice and validity of a particular analysis method. Stevens (1946, 1951) distinguished among four levels of measurement. Values on a nominal (categorical) scale only label objects. As an example, the chemical name of a compound only distinguishes it from other compounds. Values on an ordinal scale rank objects (and name them) but carry no other information. A person's finishing position in a 100-meter dash provides only ordinal information. Values on an interval scale provide information about differences between objects (and rank them and name them), such as acidity of a solution measured on the pH scale or temperature in degrees centigrade. Values on a ratio scale provide the additional information of relative size, such as the mass of an object in kilograms.

The value of recognizing the scale of a response variable lies in the guidance it provides in choosing a data analysis. However, we agree with Velleman and Wilkinson's (1993) cautions about not being too rigid in using scale to choose a data analysis. Nominal or ordered categorical data typically lead to using categorical methods for data analysis. In turn, special "distribution-free" statistical methods have been developed for ordinal data (with few ties). Both interval and ratio data tend to be considered together (at some peril in guiding a choice of valid analysis) as continuous data. Typically we will classify such data as either Gaussian or not Gaussian.

Ratio scale data are necessarily nonnegative. In such scales, a value of zero indicates the absence of the property and negative values have no meaning. In practice, such data often are positively skewed and have variance increasing with the mean, especially for a sufficient wide range of conditions. Concentration of a pollutant in a river, concentration of a drug in a person's bloodstream, and volume of lava emitted by a volcano in a month might be expected to have such a variance pattern. The data often appear Gaussian after a logarithmic or similar power transformation, including square or cube root. Muller and Fetterman (Chapter 7, 2002) described the practical use of power (Box-Cox) type transformations.

6.4 WHAT REPEATED SAMPLING SCHEME WAS USED?

The distinction between the independent sampling unit (ISU) and the observational unit plays a key role in describing the sample scheme. Table 6.2 summarizes some repeated sampling schemes. The table defines rough categories, with blurry distinctions between neighbors. Although various terms in the table are used interchangeably in the literature, the definitions and distinctions made here reflect the spirit and practice in the biological and behavioral sciences. The column for "timing" might correspond to a wide variety of dimensions, other than time, such as distance from a town, location on the surface of the earth, or amount of a treatment (measured on an interval or ratio scale).

Table 6.2 Some Simple Repeated Sampling Schemes
(*i* Indicates a Particular ISU)

Design	Number of Times	Number of ISUs	Typical Timing
Cross sectional	1	N	None
Repeated measures	$p > 1$	N	Consistent
Crossover	$p > 1$	N	Consistent
Longitudinal	$p_i > 1$	N	Inconsistent
Time series	N	1	Regular

Some terms in the table deserve clarification. Consistent timing simply requires all ISUs to be evaluated at the same times, such as Monday, Tuesday, and Thursday. In contrast, inconsistent timing allows one participant to appear on Monday, Tuesday, and Thursday, while another appears on Monday, Tuesday, and Friday. Inconsistent timing may arise due to an inability to fully control the timing of data collection. Such mistimed data often arise in human clinical trials. Regular spacing requires a constant distance between times, such as measuring air temperature at a weather station once per week for a set of N consecutive weeks. The study of univariate ANOVA models (for cross-sectional designs) led to the definition of certain terms that generalize to settings with repeated observations.

Definition 6.3 (a) A *complete* design has at least one observation per treatment combination (cell).
(b) *Balanced* designs have an equal number of observations in each cell.
(c) *Exchangeable* observations may be correlated but have identical distributions and identical relationships to other observations with which they may be exchanged.

Losing one or more independent sampling units and all associated observations usually creates unbalanced or incomplete designs. Both univariate and multivariate linear models can tolerate such deviations (in “between-subject” design) and retain excellent properties. In contrast, losing only a fraction of the observations for one or more ISUs greatly complicates the task of finding accurate estimates. Furthermore, creating accurate inferences (tests and confidence intervals) usually becomes extremely difficult in the presence of missing data, especially in small to moderate samples. The size of a sample has many possible variations in the presence of repeated measures. Most often, large-sample results refer to a setting with a fixed (or finitely bounded) number of repeated observations and an arbitrarily large number of independent sampling units.

Unfortunately, many data analysts refer to any sort of repeated sampling scheme as involving repeated measures. The failure to recognize special cases has often had the regrettable effect of more general and less accurate statistical methods being used when more accurate and easier to use methods were available.

An important special case of repeated-measures designs arises when the observations for each ISU are exchangeable; the order is arbitrary and the variable distinguishing observations within an ISU only labels them, thus providing only a nominal (categorical) scale. The term split-plot design reflects the origin of the term in agricultural statistics. A study aimed at finding the best amount of fertilizer to apply to corn might use such a design. A set of N fields from different farms represent the ISUs. With three levels of fertilizer of interest, the great variability between fields (plots of land) militates toward splitting each plot into three subplots and randomly assigning one of the three fertilizer levels to each. Local differences in such things as the quality of soil, rainfall, and cultivation equipment all contribute to between-plot variation but not within. Assuming equal variability and correlations (compound symmetry of the covariance matrix) of yields across subplots (the repeated measures) seems completely reasonable. Compound symmetry arises naturally with exchangeable observations, while nonexchangeable repeated measures rarely achieve compound symmetry.

Time as the repeated-measure dimension provides the most common example of nonexchangeable repeated observations that seem extremely unlikely to have compound symmetry. Concern about the assumption led to documentation of dramatic inflation of type I error rates under violation of the assumption of compound symmetry (Box, 1954a, b). Subsequently, statisticians developed tests robust to violation of the assumption, which implicitly allow for an arbitrary covariance structure (Geisser and Greenhouse, 1958; Greenhouse and Geisser, 1959; Huynh and Feldt, 1976) within the traditional split-plot or “univariate” approach to repeated measures. The “multivariate” test statistics, developed between roughly 1930 and 1970, avoid compound symmetry and begin with an assumption of unstructured covariance matrix.

Every analysis of repeated measures makes an implicit or explicit choice of model for the covariance pattern within a participant (ISU) across repeated measures. With Gaussian data, the validity of a choice among a univariate, multivariate, or mixed linear model depends almost entirely on the actual covariance pattern among observations. Covariance models discussed here are usually one of four types: complete independence, compound symmetry, partially structured, and completely unstructured. The four are ordered from simplest to most complex.

6.5 ANALYSIS STRATEGIES FOR MULTIVARIATE DATA

Data suitable for analysis with general linear multivariate models may be grouped into four classes, depending on scientific considerations: pure multivariate, commensurate multivariate, repeated measures (also commensurate),

and doubly multivariate. The scientific goals drive the choice of analysis strategy for any particular set of data. We focus on the following possible strategies: (1) a collection of univariate analyses with a Bonferroni correction, (2) multivariate analysis of Variance (MANOVA), (3) a collection of multivariate analyses with a Bonferroni correction, (4) repeated measures with a “multivariate” approach (MULTIREP), (5) repeated measures with a “univariate” approach (UNIREP), and (6) repeated measures with a “mixed” model approach. (7) Special purpose methods include growth curves, seemingly unrelated regression (SUR), a doubly-multivariate model (DMM), and missing data methods. Given a particular analysis strategy, more than one approach to inference may apply. Table 6.3 summarizes a number of cases.

Table 6.3 Linear Model Analysis Strategies for Gaussian Repeated Measures

Pattern of Responses	Strategy	Nominal Size of Test
p Distinct	Bonferroni univariate	α/p
	MANOVA	α
	Bonferroni MANOVA c clusters	α/c
p Repeated	Bonferroni univariate	α/p
	MULTIREP	α
	Growth curve	α
	UNIREP	α
	Mixed	α
	Bonferroni for clusters	α/c
p_1 Distinct, repeated p_2	Bonferroni 4–7 for distinct clusters	α/p_2
	Bonferroni MANOVA time clusters	α/p_1
	DMM specific method	α

A collection of univariate analyses with the total test size controlled with a Bonferroni correction may have particular appeal for a modest number of response noncommensurate variables. Additionally, no interest in profiles of response (weighted linear combinations of responses) goes with the desire to consider each variable separately. The approach allows any pattern of missing data and any variety of design matrices across responses. The allocation of test size should reflect the relative scientific interest and importance of the variables.

A collection of multivariate, MULTIREP, UNIREP, or mixed model analyses with total test size controlled by a Bonferroni correction may have particular appeal in the doubly multivariate setting. Toxicologists often study a suite of responses measured repeatedly over time. A separate multivariate could be conducted at each time. Alternately, and most commonly, a separate repeated-measures analysis could be conducted for each type of response variable.

Alternately some methods have been especially developed for the setting. Timm (2002, Section 6.7) provided a good review.

Seemingly unrelated regression (Srivastava and Giles, 1987) and growth curve models (Kshirsagar and Smith, 1995) provide generalizations of the multivariate model. SUR seeks to take advantage of correlations among responses, while allowing for different design matrices. Growth curve models seek to capitalize on modifying a multivariate approach to repeated measures to take advantage of the simplest model that fits the repeated dimension. Timm's book (2002, Chapter 5) contains a detailed introduction to both.

Table 6.4 Linear Model Form Properties for Gaussian Repeated Measures

Approach	Small- N	Σ	Additional	
	Inference?	Robust?	Plus	Minus
Bonferroni univariate	Good	Yes	Flexible	No Profile, Trend
Bonferroni multivariate	Good	Yes	Flexible	
MANOVA	Good	Yes ¹	Profiles	X Same For All No Missing
MULTIREP	Good	Yes	Trends	X Same For All No Missing
Growth curve	Good	Yes	Trends	X Same For All No Missing
UNIREP	Good ¹	Yes ¹	Trends	X Same No Missing
Mixed	Can be bad	No	Time Vary Missing OK Model Σ	X OK Fragile Limited Inference Limited Diagnostics
EM and adjusted MANOVA, MULTIREP, or UNIREP	Good	Yes	Missing OK In Small N	X Same For All
SUR	Uncertain	Yes	X Varies	No Missing Inference?
DMM	Uncertain	Uncertain	Tailored	No Missing

¹With appropriate test choice.

Table 6.4 summarizes performance characteristics of the various forms of linear models that might be used. Barton and Cramer (1989) and Catellier and Muller (2000) recommended using the EM algorithm for estimation and adjusted degree of freedom tests for MANOVA, MULTIREP, and UNIREP tests with missing data. In contrast to a mixed model, the approach always controls test size, even in very small samples. The Appendix (Section A.2) contains a description of free SAS/IML[®] code which implements the methods and where to retrieve it from the Web.

6.6 CAUTIONS AND RECOMMENDATIONS

The flexibility and generality of the mixed model make it extremely tempting to simply always use it for any linear model analysis. However, just as a skilled carpenter understands and uses many different saws, a skilled data analyst understands and uses many different kinds of linear models. Littell (2003) urged readers to recognize the limitations of the mixed model, especially in terms of accuracy of inference. Some limitations arise from computational difficulties. Numerical problems with currently popular software may badly mislead the user.

The algorithm may fail to converge to a solution, even though a valid and unique answer exists. Muller, Edwards, Simpson, and Taylor (2006) used popular mixed model software to analyze simulated Gaussian data. The observations followed a multivariate linear model with two within-subject factors, each with three levels, giving $p = 9$, and no between-subject factors. In all cases no missing or mistimed data were present and $N \in \{10, 20, 40\}$. Consequently, \mathbf{Y} was always $N \times 9$ and $\mathbf{X} = \mathbf{1}_N$ (obviously full rank). Any standard multivariate linear model program can compute the unique maximum likelihood and REML estimates for the primary parameters (\mathbf{B} and Σ), which are guaranteed to exist, in one step. For mixed model analysis, an unstructured covariance model was always requested, which ensured that the model was valid in the population. Using the default options, the program failed to converge for roughly 2% of the samples. A standard multivariate linear model program was applied to each problematic sample to verify that estimates could be computed. Merely increasing the number of iterations reduced the number of convergence failures, but not completely. Artful tuning of the convergence criteria eliminated some more of the convergence failures, but not all.

Faced with convergence failure, many data analysts would change the request for an unstructured covariance matrix to a request for a compound symmetric matrix. Although doing so might lead to convergence, the strategy will usually inflate test size in small samples. The uncorrected UNIREP test (which assumes compound symmetry) also inflates test size in the same setting but performs better than the mixed model tests in many ways when applicable.

We recommend the following steps to reduce or completely avoid such problems. (1) Use a MULTIREP or UNIREP test and associated model whenever they apply. Current mixed model tests are never better in controlling accuracy of

inference. (2) Round the values of time and related predictor values to the smallest number of digits that are scientifically meaningful. In a clinical trial with visits once per month, recording time values in days may greatly destabilize the calculations with no scientific return. Obviously the choice of recording precision must be made jointly with the scientists leading the project. (3) Given the need for a mixed model analysis, begin the process by creating well-conditioned data. Such data have (a) careful scaling, (b) removal of any location differences, (c) full-rank coding schemes for indicator variables, especially effect coding or cell mean style coding, (d) centered or pseudocentered continuous predictors, and (e) design matrices transformed to make them as close to completely orthonormal as needed. A pseudocentered variable has had a scientifically meaningful and convenient value subtracted in order to make the mean approximately zero, such as $D = T - 37$, for T human body temperature in degrees centigrade. Orthogonal or orthonormal polynomial coding for time usually helps greatly when applicable. Both X_s and Z_s , as well as y_s , should receive the improvements. Chapters 8 and 9 in Muller and Fetterman (2002) contain further discussion in the context of a univariate linear model. (4) Take advantage of the options of the particular program in use to help the program find a solution. An artful choice of starting value estimates is likely to have the most impact. The choice of algorithm and convergence criterion also can greatly affect the ability to find the solution (if one exists). (5) Finally, do not declare the covariance model to be invalid if the limitations of the data disallow estimating it. An alternate analysis method may be the only defensible choice.

Example 6.1 Analyzing data from an observational study with a univariate multiple regression model illustrates the last recommendation. An epidemiologist seeking to build a model of lung function might choose to use smoking status (with three levels, current, previous, never), race (two levels), and gender (two levels) as basic predictors. Obviously interaction variables also have appeal. However, even in a relatively large sample, it might happen that only one white female who never smoked happens to be included. Including the three-way interaction of race by gender by smoking status creates a model with severe collinearity with the intercept and extremely unstable computations. The data do not allow estimating or testing the interaction, a fact revealed by careful attention to regression diagnostics. The data do *not* provide any evidence whatsoever either in favor of or against the existence of such an interaction in the population (no inference is available because no estimate is available). An epidemiologist would report the desire to consider the interaction and the fact that the current study provides no information in either direction. The inability to fit the model in the sample at hand would not be interpreted as evidence against the interaction. Rather it indicates an inconvenient limitation stemming from the finite nature of the observational study design.

Example 6.2 A second example involving a clinical trial comparing three asthma medications also illustrates the last recommendation. Repeated clinic visits

and measurements of lung function are scheduled every month for six months. Baseline lung function (measured just before the start of treatment), treatment group, and the interaction of baseline with treatment provide the obvious starting model for between-subject effects. As usual for a repeated-measures study design for living organisms, a covariance pattern more complex than compound symmetry seems necessary. With 100 participants, recording visit time as number of days since baseline will likely create dozens of distinct times of observation. With 1000 participants, recording visit time as number of days since baseline will likely create over 100 distinct time values. In either case, even considering only relatively low order trends across time (linear, quadratic, and cubic), a mixed model analysis with an unstructured covariance pattern may fail to converge. Rounding time to week or half month will help some, as will careful use of orthogonal design coding.

In the example of the epidemiology study, difficulty came from an insufficient X matrix, while in the example of the clinical trial, difficulty came from an insufficient Z matrix. The inability to fit the model for the sample at hand should not be interpreted as evidence against an unstructured covariance pattern or in favor of a simpler model. Evidence for a simpler model might be available. A model assuming autocorrelation or a combination of autocorrelation and compound symmetry may converge *and* provide appropriate and well-behaved regression diagnostics. Without such positive results, some other options still remain viable. One simple approach would be to compute univariate analyses of trend scores and use a Bonferroni correction. Heterogeneity must be treated in a credible way with the approach. Alternately, separate univariate analyses for each time window (clinic visit number), again with a Bonferroni correction, may be preferred. Either approach has many unappealing features. We list them merely to illustrate that the desire to fit a mixed model, or any other model, does not automatically guarantee having enough data to support the model. Defensible inference from a set of data requires a scientifically credible model supported with diagnostic analysis, not merely a model that converges to a numerical solution.

6.7 REVIEW OF LINEAR MODEL NOTATION

Although originating much earlier, the theory and practice of linear models began to flower early in the 1900s. Early work used scalar notation, although proofs often included geometric representations and arguments. Increasing interest in more complex designs, repeated measures, and especially multivariate questions was coupled with a gradually increasing use of matrix notation. The advent and spread of electronic digital computers in the second half of the century accelerated the trend. By the end of the century, nearly all statisticians had access to powerful computers, and nearly all new linear model theory was cast in matrix notation.

Table 6.5 summarizes notation used for univariate, multivariate, and mixed linear models. The scalar equation provides a model statement for a single observation and is rarely used. Current discussion of linear mixed models most

often uses the “vector” form for one independent sampling unit (highlighted in gray). In contrast, current discussions of univariate and multivariate models most often use the forms for all observations (highlighted in gray). The disparity between the usual formulations of mixed and other linear models can be a source of confusion. Comparisons between models within a column of the table prove simpler and greatly help understanding. We chose notation for the models in order to facilitate such comparisons.

As discussed in Chapter 4, the growth curve model, as commonly used, can be interpreted as a special case of a restricted multivariate linear model. Although the GCM() would fit naturally into an expanded Table 6.5, the generalized GLM, the GGLM(), would not. The lack of fit arises from the fact that a GGLM() may not have any independent observations.

Table 6.6 summarizes covariance structures for the linear models described in Table 6.5. Comparing the structures for all observations considered together highlights the underlying similarities. In all cases, a block diagonal form occurs, with each diagonal block corresponding to an ISU. Zero values off diagonal reflect the statistical independence. The univariate model has all scalar diagonal blocks, all equal to σ^2 (homogeneity of variance hold). The multivariate model has all diagonal blocks of the same size and value, Σ (homogeneity of covariance holds). The mixed model allows the diagonal blocks to vary in size and value. The covariance elements must be a function of a modest number of parameters (relative to the number of observations) in order to allow computing valid estimates.

Table 6.5 Linear Model Statements

$$(Np \leftrightarrow n = \sum_{i=1}^N p_i)$$

Model	One Observation	Observations for One ISU	All Observations
$GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$	$y_i = \sum_{j=1}^q x_{i,j}\beta_j + e_i$	$y_i = \mathbf{X}_i\boldsymbol{\beta} + e_i$ $1 \times 1 \quad 1 \times q \times 1 \quad 1 \times 1$	$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ $N \times 1 \quad N \times q \times 1 \quad N \times 1$
$GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \boldsymbol{\Sigma})$	$y_{i,k} = \sum_{j=1}^q x_{i,j}\beta_{j,k} + e_{i,k}$	$\mathbf{Y}_i = \mathbf{X}_i\mathbf{B} + \mathbf{E}_i$ $1 \times p \quad 1 \times q \times p \quad 1 \times p$	$\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$ $N \times p \quad N \times q \times p \quad N \times p$
$LMM_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta}, \boldsymbol{\Sigma}_i(\boldsymbol{\tau})]$	$y_{i,k} = \sum_{j=1}^q x_{i,j}\beta_{j,k} + \sum_{l=1}^m z_{i,l}d_l + e_{i,k}$	$\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i$ $p_i \times 1 \quad p_i \times q \times 1 \quad p_i \times m \times 1 \quad p_i \times 1$	$\mathbf{y}_s = \mathbf{X}_s\boldsymbol{\beta} + \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s$ $n \times 1 \quad n \times q \times 1 \quad n \times m \times 1 \quad n \times 1$

Table 6.6 Covariance Structures for Linear Models

Model	One Observation	Observations for One ISU	All Observations
$M_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$	$\mathcal{V}(y_i) = \sigma^2$	$\mathcal{V}(y_i) = \sigma^2$	$\mathcal{V}(\mathbf{y}) = \mathbf{I}_N \sigma^2 = \mathbf{I}_N$
$M_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$	$\mathcal{V}(y_{i,k}) = \sigma_{kk}$	$\mathcal{V}(\mathbf{Y}'_i) = \boldsymbol{\Sigma}$	$\mathcal{V}[\text{vec}(\mathbf{Y}')] = \mathbf{I}_N$
$M_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i \boldsymbol{\beta}, \boldsymbol{\Sigma}_i(\boldsymbol{\tau})]$	$\mathcal{V}(y_{i,k}) = \langle \boldsymbol{\Sigma}_i(\boldsymbol{\tau}) \rangle_{kk}$	$\mathcal{V}(\mathbf{y}_i) = \boldsymbol{\Sigma}_i(\boldsymbol{\tau})$ $= \mathbf{Z}_i \boldsymbol{\Sigma}_{di}(\boldsymbol{\tau}_d) \mathbf{Z}'_i + \boldsymbol{\Sigma}_{ei}(\boldsymbol{\tau}_e)$	$\mathcal{V}(\mathbf{y}_s) = \boldsymbol{\Sigma}_s = \bigoplus_{i=1}^N$

CHAPTER 7

General Theory of Multivariate Distributions

7.1 MOTIVATION

Chapter 7 serves two purposes. First, we present a consistent notation and many basic results for groups of random variables, especially when arranged into vectors. Second, the presentation provides an implicit review of such concepts as moments and generating functions for single random variables. The reader may access the implicit review in a very simple way: Reduce all dimensions to make any matrix or vector a scalar (1×1). Applying the reduction to every result will immediately tell the reader which results are old friends that generalize conveniently and which are new acquaintances requiring extra time to get to know.

Although some results about random vectors apply directly to random matrices, many others do not. If matrix \mathbf{X} has random elements, considering $\mathbf{y} = \text{vec}(\mathbf{X})$ often proves fruitful. For symmetric \mathbf{Z} , often the form $\mathbf{u} = \text{vech}(\mathbf{Z})$ is easier to work with. In particular, $\mathbf{u} = \text{vech}(\mathbf{Z})$ may have a density, while $\text{vec}(\mathbf{Z})$ definitely does not. Gupta and Nagar's (2000) book reflects an increasing interest in studying random matrices directly. Johnson and Kotz (1972), Johnson, Kotz, and Balakrishnan (1997), Kotz, Balakrishnan, and Johnson (2000), and Kotz and Nadarajah (2004) contain more traditional treatments of joint distributions.

In the theory of distributions, the adjective “multivariate” describes any collection of more than one random variable, including ones arrayed as vectors or matrices. Clarity requires distinguishing between vector and matrix forms. An example occurs in Chapter 8, which has separate treatments of the vector Gaussian (often called the multivariate Gaussian) and the matrix Gaussian distributions.

As a general principle for the entire book, we provide a carefully stated and technically correct presentation. We rule out most pathological cases by imposing mild regularity assumptions (such as finite variance, $\sigma^2 < \infty$). We do allow for pathologies that occur naturally in practice (such as $\hat{\sigma}^2 = 0$). Taking advantage of the generality and convenience of characteristic functions involves limited consideration of complex variables. All other contexts involve only real values.

For the sake of brevity, we omit many proofs in the present chapter, especially complicated ones. Most readers will have seen scalar versions of the results in the

present chapter. The vector and matrix generalizations are not primarily linear model properties. In all subsequent chapters we (1) provide a proof, (2) provide a reference to a proof, or (3) omit the proofs with the hope that a conscientious student will derive the result.

7.2 NOTATION AND CONCEPTS

The presentation here presumes knowledge of the basic theory of probability and inference. We focus here on multivariate properties of well-defined random variables that have a joint distribution.

Most discussions consider vectors of random variables (random vectors), such as $\mathbf{y} = [y_1 \ y_2 \ \cdots \ y_n]'$, have a joint distribution defined for all $\mathbf{y}_* \in S$ (the sample space). Less often, we consider matrices of random variables (random matrices), such as $\mathbf{Y} = \{y_{jk}\}$ with a joint distribution defined for all $\mathbf{Y}_* \in S$ (the sample space). Implicitly, any result for a random vector always also applies to $\text{vec}(\mathbf{Y})$ or to $\text{vech}(\mathbf{Y})$ if $\mathbf{Y} = \mathbf{Y}'$. We will review ways of characterizing (uniquely and completely specifying) the distribution of \mathbf{y} . We will also review some general properties of distributions, such as moments. Although many results apply to any random variable, continuous random variables will be discussed in detail. In contrast, corresponding forms needed for discrete variables will mostly be omitted.

A distribution is an entity which exists apart from any particular characterization. The cumulative distribution function (CDF) is not the distribution. A well-defined distribution always has a CDF and may or may not have a probability density function (PDF). For a given CDF $F_{\mathbf{y}}(\mathbf{y}_*)$ a corresponding characteristic function (CF) $\phi_{\mathbf{y}}(\mathbf{t})$ always exists, and is unique (a.e., that is, *almost everywhere*). Conversely, if $\phi_{\mathbf{y}}(\mathbf{t})$ is the CF of a distribution, a corresponding CDF $F_{\mathbf{y}}(\mathbf{y}_*)$ always exists and is unique (a.e.). Moreover, mathematical methods allow determining one from the other. Thus a distribution is always completely and uniquely specified (a.e.) by (1) its CDF and (2) its CF.

7.3 FAMILIES OF DISTRIBUTIONS

We often consider a collection of distributions with CDFs which differ, functionally, only in the values of one or more parameters, the elements of $\boldsymbol{\theta} \in \Theta$. In particular, the univariate Gaussian family of distributions is denoted $\{\mathcal{N}(\mu, \sigma^2) : \mu^2 < \infty, 0 \leq \sigma^2 < \infty\}$, while the multivariate family is denoted $\{\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma}) : \boldsymbol{\mu} \in \mathbb{R}^p, \boldsymbol{\mu}'\boldsymbol{\mu} < \infty, (p \times p) \boldsymbol{\Sigma} = \mathbf{T}\text{Dg}(\boldsymbol{\lambda})\mathbf{T}', \lambda_j \geq 0\}$. More general families also exist. Members of the *exponential family* have PDFs which are not all of the same functional form but which can all be expressed in the general form:

$$f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) = a(\boldsymbol{\theta})b(\mathbf{y}_*)\exp\left[\sum_{j=1}^p c_j(\boldsymbol{\theta})t_j(\mathbf{y}_*)\right]. \quad (7.1)$$

Members of the *elliptically symmetric family* have PDFs depending on \mathbf{y}_* only

through a quadratic form:

$$f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) = h[(\mathbf{y}_* - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_* - \boldsymbol{\mu}); \boldsymbol{\theta}]. \quad (7.2)$$

Members of the *spherically symmetric family* have PDFs depending on \mathbf{y} only through an inner product:

$$f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) = h(\mathbf{y}_*' \mathbf{y}_*; \boldsymbol{\theta}) = h[(\mathbf{y}_* - \mathbf{0})' \mathbf{I} (\mathbf{y}_* - \mathbf{0}); \boldsymbol{\theta}]. \quad (7.3)$$

7.4 CUMULATIVE DISTRIBUTION FUNCTION

Definition 7.1 Every distribution is characterized by a *joint cumulative distribution function* (CDF):

$$\begin{aligned} F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) &= \Pr\{\mathbf{y} \leq \mathbf{y}_* | \boldsymbol{\theta}\} \\ &= \Pr\left\{ \bigcap_{j=1}^n (y_j \leq y_{j*}) | \boldsymbol{\theta} \right\}. \end{aligned} \quad (7.4)$$

The simplified notation $F_{\mathbf{y}}(\mathbf{y}_*)$ is also used and the CDF is also known as the *distribution function*.

Definition 7.2 A random vector \mathbf{y} will be described as *discrete* if and only if it has countable support (smallest S such that $\Pr\{\mathbf{y} \in S\} = 1$).

Definition 7.3 $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ and \mathbf{y} are *continuous* if and only if $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ is a continuous function of \mathbf{y}_* .

7.5 PROBABILITY DENSITY FUNCTION

Definition 7.4 $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ and \mathbf{y} are *absolutely continuous* (a.c.) if and only if a nonnegative function $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ exists such that the following n -fold integral over region $S_{\mathbf{y}_0} = \{\mathbf{y}_* : \mathbf{y}_* \leq \mathbf{y}_0\}$ exists $\forall \mathbf{y}_0 \in \mathbb{R}^n$:

$$F_{\mathbf{y}}(\mathbf{y}_0; \boldsymbol{\theta}) \equiv \int_{S_{\mathbf{y}_0}} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_*. \quad (7.5)$$

The $d\mathbf{y}_*$ means $dy_{*1} dy_{*2} \cdots dy_{*n}$ in the n -fold integral.

Definition 7.5 For continuous distributions, the partial derivative

$$f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) = \frac{\partial^{(n)}}{\partial y_{*1} \partial y_{*2} \cdots \partial y_{*n}} F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) \quad (7.6)$$

is the *probability density function* (PDF), or *the density* of the distribution, if it exists $\forall \mathbf{y}_* \in S$.

Theorem 7.1 If the PDF $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ exists then, with $S_{\mathbf{y}_0} = \{\mathbf{y}_* : \mathbf{y}_* \leq \mathbf{y}_0\}$,

1. $\int_{\mathbb{R}^n} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_* = 1$ (7.7)
2. $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) \geq 0 \forall \mathbf{y}_* \in \mathbb{R}^n$
3. $F_{\mathbf{y}}(\mathbf{y}_0; \boldsymbol{\theta}) = \int_{S_{\mathbf{y}_0}} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_*$.

Any function satisfying 1 and 2 is a PDF. Condition 3 implies it is a PDF of $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$.

Theorem 7.2 If $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ is a PDF, and function $g(\mathbf{y}_*; \boldsymbol{\theta}) = f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ except at a countable number of points, then $g(\mathbf{y}_*; \boldsymbol{\theta})$ is also a PDF of $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$. Thus we say $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ is unique (a.e.).

The nonuniqueness must be considered in defining maximum likelihood estimation.

7.6 FORMULAS FOR PROBABILITIES AND MOMENTS

The formula for a probability or moment changes with the type of random variable. Discrete random variables allow writing the formula in terms of a (possibly infinite) summation, while absolutely continuous random variables use an integral. However, some random variables, such as censored survival time of a person in a clinical trial, are neither discrete nor continuous. Furthermore, some continuous distributions do not have a PDF (although the CDF always exists).

A straightforward generalization of the usual (Riemann) integral solves the problem for all random variables of interest in the present book. The Stieltjes (or Riemann-Stieltjes) integral allows defining probability and moment formulas for nearly any type of random variable (discrete, absolutely continuous, or “neither”). Handling the mathematically exotic random variables not covered by the Stieltjes integral requires measure theory and the methods of Lebesgue integration.

Definition 7.6 Following Weisstein (2003), a Stieltjes integral applies to real functions $f(x)$ and $h(x)$, both bounded on the closed interval $[a, b]$. For a partition $\{a = x_0 < x_1 < \dots < x_{n-1} < x_n = b\}$ and $x_j < e_j < x_{j+1}$,

$$s = \sum_{j=0}^{n-1} f(e_j)[h(x_{j+1}) - h(x_j)] \tag{7.8}$$

is a Riemann sum. If $s \rightarrow s_0$, a fixed number, as $\max(x_{j+1} - x_j) \rightarrow 0$, then s_0 is the *Stieltjes integral*, written

$$s_0 = \int_a^b f(x) dh(x). \tag{7.9}$$

Considering the complex plane extends the integral to complex variables.

Continuity of $f(x)$ and bounded variation of $h(x)$ over $[a, b]$ ensure the Stieltjes integral exists for $[a, b]$. It fails to exist if $f(x)$ and $h(x)$ are not continuous at a common point. If $h(x)$ has a continuous derivative, the Stieltjes integral reduces to the Riemann integral. Any finite or infinite sum can be expressed as a Stieltjes integral with an appropriate choice of $h(x)$.

7.7 CHARACTERISTIC FUNCTION

Definition 7.7 The *characteristic function* (CF) of the distribution of \mathbf{y} is defined for $i = \sqrt{-1}$ and $\forall \mathbf{t} \in \mathbb{R}^n$ as

$$\begin{aligned} \phi_{\mathbf{y}}(\mathbf{t}) &= E[\exp(i\mathbf{t}'\mathbf{y})] = E[\cos(\mathbf{t}'\mathbf{y})] + iE[\sin(\mathbf{t}'\mathbf{y})] \\ &= \int_{\mathbb{R}^n} \exp(i\mathbf{t}'\mathbf{y}_*) dF_{\mathbf{y}}(\mathbf{y}_*) \\ &= \int_{\mathbb{R}^n} \exp(i\mathbf{t}'\mathbf{y}_*) f_{\mathbf{y}}(\mathbf{y}_*) d\mathbf{y}_*, \end{aligned} \tag{7.10}$$

with the last line applying only if $f_{\mathbf{y}}(\mathbf{y}_*)$ exists. The first two always apply if interpreted in terms of complex-variable Stieltjes integration.

In general the integrals in the definition require contour integration methods (taught in a class on complex variables). The last integral is the Fourier transform of $f_{\mathbf{y}}(\mathbf{y}_*)$. Here $e^{iz} \equiv \cos(z) + i \cdot \sin(z)$ is a complex number. It is represented in two dimensions by the coordinates of a point located on the unit circle, $[\cos(z), \sin(z)]$. Its magnitude is $|e^{iz}| = \sqrt{\cos^2(z) + \sin^2(z)} \leq 1$. The mapping of $[\cos(\mathbf{t}'\mathbf{y}), \sin(\mathbf{t}'\mathbf{y})]$ has a simple geometric interpretation (Epps, 1993). The mapping wraps the PDF around the unit circle, while $\phi_{\mathbf{y}}(\mathbf{t}) = E(\cdot)$ specifies the location of the center of mass within the circle.

The definition makes it obvious that derivation and manipulation of characteristic functions require complex-variable analysis. Despite that, characteristic functions can provide many powerful and practical results with knowledge of only the most basic results about the complex variables. Most importantly for the developments in later chapters, linear transformations of random variables induce simple changes in characteristic functions of random vectors and matrices. With Gaussian errors, the use of characteristic functions greatly simplifies the derivation of distributions of mean and covariance estimators. The approach has the important advantage of allowing less-than-full-rank coding schemes and population covariance matrices with little extra work.

Theorem 7.3 (a) For any random vector \mathbf{y} the characteristic function $\phi_{\mathbf{y}}(\mathbf{t})$ always exists finitely; specifically, $|\phi_{\mathbf{y}}(\mathbf{t})| \leq 1 \forall \mathbf{t} \in \mathbb{R}^n$.

(b) Any multivariate characteristic function $\phi_{\mathbf{y}}(\mathbf{t})$ is a uniformly continuous function. Equivalently, $\lim_{\mathbf{h} \rightarrow 0} |\phi_{\mathbf{y}}(\mathbf{t} + \mathbf{h}) - \phi_{\mathbf{y}}(\mathbf{t})| = 0$.

(c) If \mathbf{y} ($n \times 1$) is distributed with CDF $F_{\mathbf{y}}(\mathbf{y}_*)$, characteristic function $\phi_{\mathbf{y}}(\mathbf{t})$, and $F_{t\mathbf{y}}(s)$ specifies a distribution symmetric about zero $\forall \mathbf{t} \in \mathbb{R}^n$, then the multivariate characteristic function is real valued and $\phi_{\mathbf{y}}(\mathbf{t}) \in [-1, 1]$.

Proof of (a). If $\mathbf{t} = \mathbf{0}$, then $\phi_{\mathbf{y}}(\mathbf{0}) = E(e^{i \cdot 0}) = E(1) = 1$. If $\mathbf{t} \neq \mathbf{0}$, we consider a fixed value of \mathbf{t} . Since \mathbf{t} is fixed, $a \in \mathbb{R}$ exists such that $\phi_{\mathbf{y}}(\mathbf{t}) = |\phi_{\mathbf{y}}(\mathbf{t})|e^{ia}$ (the "polar form"). Constant a is an unknown function of \mathbf{t} . For fixed $\{\mathbf{t}, a\}$, $\exists \mathbf{c}$ such that $a = \mathbf{t}'\mathbf{c}$. The choice $\mathbf{c} = \mathbf{t}(a/\mathbf{t}'\mathbf{t})$ will do. It follows that $\phi_{\mathbf{y}}(\mathbf{t}) = |\phi_{\mathbf{y}}(\mathbf{t})|e^{i\mathbf{t}'\mathbf{c}}$ and $|\phi_{\mathbf{y}}(\mathbf{t})| = \phi_{\mathbf{y}}(\mathbf{t})e^{-i\mathbf{t}'\mathbf{c}}$. Thus $|\phi_{\mathbf{y}}(\mathbf{t})| = E[\exp(i\mathbf{t}'\mathbf{y} - i\mathbf{t}'\mathbf{c})] = \phi_{(\mathbf{y}-\mathbf{c})}(\mathbf{t})$. Since $|\phi_{\mathbf{y}}(\mathbf{t})| \in \mathbb{R}$ by definition, $\phi_{(\mathbf{y}-\mathbf{c})}(\mathbf{t}) = E[\cos[\mathbf{t}'(\mathbf{y} - \mathbf{c})] + i \cdot 0]$. Necessarily $|\phi_{\mathbf{y}}(\mathbf{t})| \in [-1, 1]$ since $\cos(\cdot) \in [-1, 1]$.

Proof of (b) is left as an exercise.

Proof of (c). It is left as an exercise to prove the result holds for $n = 1$. Hence it holds for univariate characteristic functions in the set $\{\phi_{y_j}(s) : s \in \mathbb{R}, \forall y_j\}$, and for all characteristic functions in the set $\{\phi_{t\mathbf{y}}(s) : \mathbf{t} \in \mathbb{R}^n, s \in \mathbb{R}\}$. Necessarily $\phi_{t\mathbf{y}}(s) \in [-1, 1] \forall \mathbf{t} \in \mathbb{R}^n, s \in \mathbb{R}$, including $s = 1$. Also, $\phi_{\mathbf{y}}(\mathbf{t}) \in [-1, 1] \forall \mathbf{t} \in \mathbb{R}^n$ since $\phi_{t\mathbf{y}}(1) = \phi_{\mathbf{y}}(\mathbf{t}) \forall \mathbf{t} \in \mathbb{R}^n$. \square

Many properties of characteristic functions of multivariate distributions can be proven via the following lemma and theorems. In the following discussions, sometimes it will be helpful to have $\phi(s; \mathbf{t})$ denote the characteristic function of the univariate random variable $\mathbf{t}'\mathbf{y}$ and write $\phi(s; \mathbf{t}) = \phi_{t\mathbf{y}}(s) = E\{\exp[is(\mathbf{t}'\mathbf{y})]\}$. The notation emphasizes that s is the argument of function ϕ while \mathbf{t} is interpreted as a fixed parameter.

Lemma 7.1 If $n \times 1$ \mathbf{y} has a multivariate distribution and \mathbf{t} is a nonrandom $n \times 1$ vector $\in \mathfrak{R}^n$, then $z = \mathbf{t}'\mathbf{y} = \sum_{j=1}^n t_j y_j$ is a scalar random variable with a univariate distribution.

Theorem 7.4 Knowing $\{\phi_{\mathbf{t}'\mathbf{y}}(s) : \mathbf{t} \in \mathfrak{R}^n, s = 1\}$ determines $\phi_{\mathbf{y}}(\mathbf{t})$ (Cramér-Wold).

Proof. By definition, the characteristic function of random variable $\mathbf{t}'\mathbf{y}$ is $\phi_{\mathbf{t}'\mathbf{y}}(s) = E\{\exp[is(\mathbf{t}'\mathbf{y})]\} \forall s \in \mathfrak{R}$ and the characteristic function of random vector \mathbf{y} is $\phi_{\mathbf{y}}(\mathbf{t}) = E[\exp(i\mathbf{t}'\mathbf{y})] \forall \mathbf{t} \in \mathfrak{R}^n$. By evaluating the univariate characteristic function at $s = 1$ we can determine the value of the multivariate characteristic function $\forall \mathbf{t} \in \mathfrak{R}^n$. In fact, $\phi_{\mathbf{t}'\mathbf{y}}(1) = \phi_{\mathbf{y}}(\mathbf{t}) \forall \mathbf{t} \in \mathfrak{R}^n$. \square

Theorem 7.5 Knowing $\phi_{\mathbf{y}}(\mathbf{t})$ determines $\{\phi_{\mathbf{t}_0'\mathbf{y}}(s) : \mathbf{t}_0 \in \mathfrak{R}^n, s \in \mathfrak{R}\}$.

Proof. Here \mathbf{t}_0 is any fixed, arbitrary vector in \mathfrak{R}^n . By definition the characteristic function of random variable $\mathbf{t}_0'\mathbf{y}$ is $\phi_{\mathbf{t}_0'\mathbf{y}}(s) = E\{\exp[is(\mathbf{t}_0'\mathbf{y})]\} \forall s \in \mathfrak{R}$ and the characteristic function of random vector \mathbf{y} is $\phi_{\mathbf{y}}(\mathbf{t}) = E[\exp(i\mathbf{t}'\mathbf{y})] \forall \mathbf{t} \in \mathfrak{R}^n$. By evaluating the multivariate characteristic function at $\mathbf{t} = s\mathbf{t}_0$ we can determine the value of the univariate characteristic function $\forall s \in \mathfrak{R}$. Merely evaluate $\phi_{\mathbf{y}}(s\mathbf{t}_0) = \phi_{\mathbf{t}_0'\mathbf{y}}(s) \forall \mathbf{t}_0 \in \mathfrak{R}^n$. \square

Theorem 7.6 If the distribution of \mathbf{y} ($n \times 1$) is absolutely continuous, then the PDF is determined by the characteristic function as

$$f_{\mathbf{y}}(\mathbf{y}_*) = (2\pi)^{-n} \int_{\mathfrak{R}^n} \exp(-i\mathbf{t}'\mathbf{y}_*) \phi_{\mathbf{y}}(\mathbf{t}) dt. \quad (7.11)$$

Proof. The proof is left as an exercise.

Theorem 7.7 If x is distributed with CDF $F_x(x_*)$ and characteristic function $\phi_x(t)$ while a and b are two points of continuity of $F_x(x_*)$, then $\Pr\{a \leq x \leq b\} \equiv F_x(b) - F_x(a)$ can be expressed as

$$\Pr\{a \leq x \leq b\} = \lim_{T \rightarrow \infty} \frac{1}{2\pi} \int_{-T}^T \frac{\exp(-ita) - \exp(-itb)}{-it} \phi_x(t) dt. \quad (7.12)$$

The result implies $\phi_x(s)$ completely determines $F_x(x_*)$.

Proof. The proof is left as an exercise.

The form in the last theorem does not directly provide the CDF if $\Pr\{x < 0\} > 0$ and $\Pr\{x > 0\} > 0$. Imhof (1961) recommended the following result for such cases.

Theorem 7.8 If random variable x has CDF $F_x(x_*)$ and characteristic function $\phi_x(t)$, while $\Im(z)$ indicates the imaginary part of z , then for x_0 a point of continuity of $F_x(x_0)$,

$$F_x(x_0) = \lim_{T \rightarrow \infty} \frac{1}{\pi} \int_0^T \frac{\Im[\exp(-itx_0)\phi_x(t)]}{t} dt. \quad (7.13)$$

Proof. Gil-Pelaez (1951) included a proof.

Theorem 7.9 If \mathbf{y} ($n \times 1$) is distributed with CDF $F_{\mathbf{y}}(\mathbf{y}_*)$, characteristic function $\phi_{\mathbf{y}}(\mathbf{t})$, and $\int_{\mathbb{R}^n} |\phi_{\mathbf{y}}(\mathbf{t})| d\mathbf{t} < \infty$, then $F_{\mathbf{y}}(\mathbf{y}_*)$ is absolutely continuous and a PDF $f_{\mathbf{y}}(\mathbf{y}_*)$ exists.

Proof. The proof is left as an exercise.

Theorem 7.10 If random vector \mathbf{x} has characteristic function $\phi_{\mathbf{x}}(\mathbf{t})$ and $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$, for conforming constants \mathbf{A} and \mathbf{b} , then

$$\phi_{\mathbf{y}}(\mathbf{s}) = \exp(i\mathbf{s}'\mathbf{b})\phi_{\mathbf{x}}(\mathbf{A}'\mathbf{s}). \quad (7.14)$$

Proof. $\phi_{\mathbf{y}}(\mathbf{s}) = E\{\exp[i\mathbf{s}'(\mathbf{A}\mathbf{x} + \mathbf{b})]\} = \exp(i\mathbf{s}'\mathbf{b})E[\exp(i\mathbf{s}'\mathbf{A}\mathbf{x})] = \exp(i\mathbf{s}'\mathbf{b})E\{\exp[i(\mathbf{A}'\mathbf{s})'\mathbf{x}]\} = \exp(i\mathbf{s}'\mathbf{b})\phi_{\mathbf{x}}(\mathbf{A}'\mathbf{s}).$ □

The notation $x \perp\!\!\!\perp y$ indicates the random variables x and y are statistically independent.

Theorem 7.11 If $\mathbf{x}_1 \perp\!\!\!\perp \mathbf{x}_2$ and $\mathbf{y} = \mathbf{x}_1 + \mathbf{x}_2$, then $\phi_{\mathbf{y}}(\mathbf{t}) = \phi_{\mathbf{x}_1}(\mathbf{t})\phi_{\mathbf{x}_2}(\mathbf{t})$.

Proof. $\phi_{\mathbf{y}}(\mathbf{t}) = E[e^{i\mathbf{t}'(\mathbf{x}_1 + \mathbf{x}_2)}] = E(e^{i\mathbf{t}'\mathbf{x}_1} e^{i\mathbf{t}'\mathbf{x}_2}) = E(e^{i\mathbf{t}'\mathbf{x}_1})E(e^{i\mathbf{t}'\mathbf{x}_2}).$ □

Theorem 7.12 If $E(|y|^m) < \infty$, then the derivative of order m of $\phi_y(t)$ exists $\forall t$ and is a uniformly continuous function with $\phi_y^{(m)}(t) \Big|_{t=0} = i^m E(y^m)$.

A converse is true for even but not odd m .

Proof. Feller (1968) provided a proof.

We have given only a brief introduction to characteristic functions. Kendall and Stuart (1977) provided a detailed presentation, including proofs of many of the results not proven here. Lukacs (1983) reviewed subsequent developments. Epps (1993) gave an excellent tutorial on the interpretation and value of characteristic functions based on geometric intuition (in the complex plane).

7.8 MOMENT GENERATING FUNCTION

Definition 7.8 If it exists for all real $t_j \in \{-t_{*j}, t_{*j}\}$ with $t_{*j} > 0$, the *moment generating function* (MGF) of the distribution of random vector $(n \times 1)$ \mathbf{y} is

$$m_{\mathbf{y}}(\mathbf{t}) = E[\exp(\mathbf{t}'\mathbf{y})]. \tag{7.15}$$

Unlike the characteristic function, the MGF does *not* exist for every distribution. However, when the MGF does exist, the characteristic function is the MGF with $i\mathbf{t}$ replacing \mathbf{t} , namely $\phi_{\mathbf{y}}(\mathbf{t}) = m_{\mathbf{y}}(i\mathbf{t})$. A proof requires consideration of the concept of analytic continuation.

Theorem 7.13 If the random vector \mathbf{y} has mean $\boldsymbol{\mu}$, dispersion $\boldsymbol{\Sigma}$, and MGF $m_{\mathbf{y}}(\mathbf{t}) = E[\exp(\mathbf{t}'\mathbf{y})]$, then

$$\begin{aligned} \left. \frac{\partial m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t}} \right|_{\mathbf{t}=\mathbf{0}} &= E(\mathbf{y}) \\ &= \boldsymbol{\mu} \end{aligned} \tag{7.16}$$

and

$$\begin{aligned} \left. \frac{\partial^2 m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}'} \right|_{\mathbf{t}=\mathbf{0}} &= E(\mathbf{y}\mathbf{y}') \\ &= \boldsymbol{\Sigma} + \boldsymbol{\mu}\boldsymbol{\mu}'. \end{aligned} \tag{7.17}$$

When they exist, $E(\mathbf{y}\mathbf{y}') = \mathcal{V}(\mathbf{y}) + E(\mathbf{y})E(\mathbf{y}')$ and $\partial e^{\boldsymbol{\mu}'\mathbf{t}}/\partial \mathbf{t} = e^{\boldsymbol{\mu}'\mathbf{t}} \partial \boldsymbol{\mu}'\mathbf{t}/\partial \mathbf{t}$.

Proof. The first derivative is an $n \times 1$ vector,

$$\begin{aligned} \frac{\partial m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t}} &= \frac{\partial}{\partial \mathbf{t}} E[\exp(\mathbf{t}'\mathbf{y})] \\ &= E \left[\frac{\partial}{\partial \mathbf{t}} \exp(\mathbf{t}'\mathbf{y}) \right] = E[\exp(\mathbf{t}'\mathbf{y})\mathbf{y}]. \end{aligned} \tag{7.18}$$

In turn, evaluating $\partial m_{\mathbf{y}}(\mathbf{t})/\partial \mathbf{t}$ at $\mathbf{t} = \mathbf{0}$ gives $E(\mathbf{1} \cdot \mathbf{y}) = E(\mathbf{y})$.

The second derivative is an $n \times n$ matrix,

$$\begin{aligned} \frac{\partial^2 m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}'} &= \frac{\partial}{\partial \mathbf{t}} \left[\frac{\partial m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t}'} \right]' = \frac{\partial}{\partial \mathbf{t}} E[\exp(\mathbf{t}'\mathbf{y})\mathbf{y}'] \\ &= E \left[\frac{\partial}{\partial \mathbf{t}} \exp[(\mathbf{t}'\mathbf{y})\mathbf{y}'] \right] \\ &= E \left[\exp(\mathbf{t}'\mathbf{y}) \frac{\partial}{\partial \mathbf{t}} \mathbf{y}' + \frac{\partial}{\partial \mathbf{t}} \exp(\mathbf{t}'\mathbf{y})\mathbf{y}' \right] \\ &= E[\exp(\mathbf{t}'\mathbf{y}) \cdot \mathbf{0} + \exp(\mathbf{t}'\mathbf{y})\mathbf{y}\mathbf{y}'], \end{aligned} \tag{7.19}$$

with the matrix \mathbf{O} being $n \times n$. In turn,

$$\begin{aligned} \left. \frac{\partial^2 m_{\mathbf{y}}(\mathbf{t})}{\partial \mathbf{t} \partial \mathbf{t}'} \right|_{\mathbf{t}=\mathbf{0}} &= \mathbf{E}(\mathbf{O} + \mathbf{1} \cdot \mathbf{y}\mathbf{y}') & (7.20) \\ &= \mathbf{E}(\mathbf{y}\mathbf{y}') \equiv \mathcal{V}(\mathbf{y}) + \mathbf{E}(\mathbf{y})\mathbf{E}(\mathbf{y}'). & \square \end{aligned}$$

Theorem 7.14 If \mathbf{x} is a random vector with MGF $m_{\mathbf{x}}(\mathbf{t})$ and $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ for conforming constants \mathbf{A} and \mathbf{b} , then

$$m_{\mathbf{y}}(\mathbf{s}) = \exp(\mathbf{s}'\mathbf{b})m_{\mathbf{x}}(\mathbf{A}'\mathbf{s}). \tag{7.21}$$

Proof. $m_{\mathbf{y}}(\mathbf{s}) = \mathbf{E}\{\exp[\mathbf{s}'(\mathbf{A}\mathbf{x} + \mathbf{b})]\} = \exp(\mathbf{s}'\mathbf{b})\mathbf{E}[\exp(\mathbf{s}'\mathbf{A}\mathbf{x})] = \exp(\mathbf{s}'\mathbf{b})\mathbf{E}[\exp(\mathbf{A}'\mathbf{s})'\mathbf{x}] = \exp(\mathbf{s}'\mathbf{b})m_{\mathbf{x}}(\mathbf{A}'\mathbf{s}).$ □

7.9 CUMULANT GENERATING FUNCTION

Definition 7.9 (a) The notation $m_{\mathbf{x}}(\mathbf{t})$ and $\phi_{\mathbf{x}}(\mathbf{t})$ indicates the MGF (if it exists) and characteristic function, respectively, of the random vector \mathbf{x} . If $m_{\mathbf{x}}(\mathbf{t})$ exists, then

$$c_{\mathbf{x}}(\mathbf{t}) = \log[m_{\mathbf{x}}(\mathbf{t})] \tag{7.22}$$

or

$$c_{\mathbf{x}}(\mathbf{t}) = \log[\phi_{\mathbf{x}}(\mathbf{t})] \tag{7.23}$$

may be defined to be the *cumulant generating function* (CGF) $c_{\mathbf{x}}(\mathbf{t})$ of the distribution (depending on the author).

(b) If a scalar random variable x has a power series expansion for $\log[m_x(t)]$, then

$$\log[m_x(t)] = \sum_{m=0}^{\infty} \kappa_m \frac{t^m}{m!}. \tag{7.24}$$

The coefficients $\{\kappa_m\}$ are the *cumulants*.

Except for the first cumulant, which equals the mean, all cumulants of the random variable x are also the cumulants of the random variable $x + a$ for constant a (hence the term *semi-invariants*). After taking derivatives, cumulant m equals the derivative of order m of $\log[m(t)]$ evaluated at $t = 0$,

$$\kappa_m = \left. \frac{\partial^{(m)} \log[m_x(t)]}{\partial t^{(m)}} \right|_{t=0}. \tag{7.25}$$

If a distribution is determined by its moments, then it is also determined by its cumulants. Not all distributions are completely determined by their moments.

Moments (and therefore moment and cumulant generating functions) fully characterize one class of random variables, namely all with finite variation. If $\Pr\{-\infty < a \leq x \leq b < \infty\} = 1$ then, $0 \leq E(|x|^m) < \infty$. For some distributions, cumulants and the CGF are easier to manipulate than moments and the MGF.

A useful application of cumulants occurs in the derivation of moments of certain distributions. A moment about zero is defined as $\mu'_m = E(x^m)$, while a central moment is $\mu_m = E\{[x - E(x)]^m\}$. Although only positive integer values of m are considered for cumulants, moments may exist (or may not, depending on the random variable) for any real m . When the moments and cumulants of order m exist, κ_m may be written as a polynomial of degree m in $\{\mu'_1, \dots, \mu'_m\}$ or $\{\mu_1, \dots, \mu_m\}$; alternately, μ_m or μ'_m may be written as a polynomial of degree m in $\{\kappa_1, \dots, \kappa_m\}$ (Kendall and Stuart, 1977, vol. 1, Section 3.1.4). Given the first four cumulants, the first four moments can easily be found because

$$\kappa_1 = \mu'_1 = E[(x - 0)^1] = \mu = E(x) \tag{7.26}$$

$$\kappa_2 = \mu_2 = E[(x - \mu)^2] = \sigma^2 = \mathcal{V}(x) \tag{7.27}$$

$$\kappa_3 = \mu_3 = E[(x - \mu)^3] \tag{7.28}$$

$$\kappa_4 = \mu_4 - 3\mu^2 = E[(x - \mu)^4] - 3\mu^2. \tag{7.29}$$

Kendall and Stuart (1977), Johnson, Kotz, and Balakrishnan (1994), and Harvey (1972) presented further details.

Theorem 7.15 If random vector \mathbf{y}_j ($n \times 1$) has CGF $c_{\mathbf{y}_j}(\mathbf{t}_j)$ and m finite cumulants for (finite) $j \in \{1, 2, \dots, J\}$, with $\{\mathbf{y}_j\}$ mutually independent (Section 7.12 contains a precise definition of “mutual independence for vectors”) and $\mathbf{s} = \sum_{j=1}^J \mathbf{y}_j$, then

$$c_{\mathbf{s}}(\mathbf{t}) = \sum_{j=1}^J c_{\mathbf{y}_j}(\mathbf{t}). \tag{7.30}$$

If $n = 1$, then

$$\kappa_m(\mathbf{s}) = \sum_{j=1}^J \kappa_m(y_j). \tag{7.31}$$

The theorem provides some extremely practical and convenient formulas. When the expressions make sense, the CGF of a sum is the sum of the component CGFs. Also, a cumulant of a sum is the sum of the corresponding cumulants. The forms have value in both analytic and numerical work.

7.10 TRANSFORMING RANDOM VARIABLES

Theorem 7.16 If \mathbf{y} ($n \times 1$) is distributed with PDF $f_{\mathbf{y}}(\mathbf{y}_*)$ and $\mathbf{t}(\mathbf{y})$ is a vector-valued function defining $(n \times 1)$ $\mathbf{z} = \mathbf{t}(\mathbf{y})$, with PDF $f_{\mathbf{z}}(\mathbf{z}_*)$ such that the mapping of points in the support of \mathbf{y} onto points in the support of \mathbf{z} is one-to-one (a.e.), then $\mathbf{z} = \mathbf{t}(\mathbf{y})$ and $\mathbf{y} = \mathbf{t}^{-1}(\mathbf{z})$ exists. Furthermore, with notation as in Searle (1982),

$$\begin{aligned} \mathbf{J}_{\mathbf{y}_* \rightarrow \mathbf{z}_*} &= \left\{ \frac{\partial y_{*j}}{\partial z_{*k}} \right\} = \{j_{jk}\} \\ &= \frac{\partial \mathbf{y}_*}{\partial \mathbf{z}'_*} = \left(\frac{\partial \mathbf{y}'_*}{\partial \mathbf{z}_*} \right)', \end{aligned} \quad (7.32)$$

the $(n \times n)$ *Jacobian* matrix for the transformation from \mathbf{y}_* to \mathbf{z}_* , exists. Indicating the absolute value of the determinant of the Jacobian as

$$J_{\mathbf{y}_* \rightarrow \mathbf{z}_*} = \|\mathbf{J}_{\mathbf{y}_* \rightarrow \mathbf{z}_*}\| = \frac{1}{\|\mathbf{J}_{\mathbf{z}_* \rightarrow \mathbf{y}_*}\|} \quad (7.33)$$

allows writing

$$f_{\mathbf{z}}(\mathbf{z}_*) = f_{\mathbf{y}}[\mathbf{t}^{-1}(\mathbf{z}_*)] J_{\mathbf{y}_* \rightarrow \mathbf{z}_*}. \quad (7.34)$$

Many authors (including Schott, 2005) define $\mathbf{J}_{\mathbf{z}_* \rightarrow \mathbf{y}_*}$ as the Jacobian. The reciprocal relationship for respective determinants ensures that either definition leads to the same density for the transformed variables.

Definition 7.10 (a) If \mathbf{z} ($n \times 1$) is a random vector and \mathbf{A} ($n \times n$) is a finite constant matrix, then $\mathbf{z} = \mathbf{A}\mathbf{y}$ is a *linear transformation*.

(b) If constant (and finite) $\mathbf{b} \neq \mathbf{0}$ is $n \times 1$ then $\mathbf{z} = \mathbf{y} + \mathbf{b}$ is a *translation* (shift in origin). Defining $\mathbf{x} = \mathbf{A}\mathbf{y} + \mathbf{b}$ indicates a linear transformation with a translation.

(c) A linear transformation is *full rank* (or *nonsingular*) when \mathbf{A} is full rank (nonsingular). In such cases \mathbf{A}^{-1} exists uniquely and the transformation is one to one and invertible. The transformation is *less than full rank* (or *singular*) when \mathbf{A} is less than full rank (singular). In such cases \mathbf{A}^{-1} does not exist and the transformation is neither one to one nor invertible.

The definitions reflect precise, mathematical descriptions of transformations. More loosely, it is often convenient to describe a linear transformation with translation as simply a linear transformation. Describing a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ or $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\boldsymbol{\beta}, \boldsymbol{\Sigma})$ as a “linear” model because the unknown parameters enter the model equation linearly is consistent with the looser usage.

Lemma 7.2 If $\mathbf{z}_* = \mathbf{A}\mathbf{y}_* + \mathbf{b}$ and \mathbf{A}^{-1} ($n \times n$) exists, then $\mathbf{y}_* = \mathbf{A}^{-1}(\mathbf{z}_* - \mathbf{b})$,

$$J_{\mathbf{y}_* \rightarrow \mathbf{z}_*} = \|\mathbf{J}_{\mathbf{y}_* \rightarrow \mathbf{z}_*}\| = \|\partial \mathbf{y}_* / \partial \mathbf{z}'_*\| = \|\mathbf{A}^{-1}\|, \quad (7.35)$$

and

$$J_{\mathbf{z}_* \rightarrow \mathbf{y}_*} = \|\mathbf{J}_{\mathbf{z}_* \rightarrow \mathbf{y}_*}\| = \|\partial \mathbf{z}_* / \partial \mathbf{y}'_*\| = \|\mathbf{A}\|. \quad (7.36)$$

Here $\|\mathbf{M}\|$ indicates the absolute value of the determinant of \mathbf{M} .

The following theorem illustrates the use of the transformation theorem for a linear transformation. The result is used frequently.

Theorem 7.17 (a) If \mathbf{y} is distributed with PDF $f_{\mathbf{y}}(\mathbf{y}_*)$, $n \times n$ constant \mathbf{A} is full rank and finite, and $n \times 1$ constant \mathbf{b} is finite, then the linear transformation $\mathbf{z} = \mathbf{A}\mathbf{y} + \mathbf{b} \equiv \mathbf{t}(\mathbf{y})$ is one to one and $\mathbf{y} = \mathbf{A}^{-1}(\mathbf{z} - \mathbf{b}) \equiv \mathbf{t}^{-1}(\mathbf{z})$.

(b) The Jacobian of the transformation from \mathbf{y} to \mathbf{z} is $J_{\mathbf{y}_* \rightarrow \mathbf{z}_*} = \|\mathbf{J}_{\mathbf{y}_* \rightarrow \mathbf{z}_*}\|$ with $\mathbf{J}_{\mathbf{y}_* \rightarrow \mathbf{z}_*} = \partial \mathbf{y} / \partial \mathbf{z}' = \mathbf{A}^{-1}$.

(c) Also,

$$\int_{\mathbb{R}^n} f_{\mathbf{y}}(\mathbf{y}_*) d\mathbf{y}_* = \int_{\mathbb{R}^n} f_{\mathbf{y}}[\mathbf{t}^{-1}(\mathbf{z}_*)] J_{\mathbf{y}_* \rightarrow \mathbf{z}_*} d\mathbf{z}_*. \quad (7.37)$$

(d) Furthermore $\mathbf{z} = \mathbf{A}\mathbf{y} + \mathbf{b}$ is distributed with PDF

$$f_{\mathbf{z}}(\mathbf{z}_*) = f_{\mathbf{y}}[\mathbf{A}^{-1}(\mathbf{z}_* - \mathbf{b})] \|\mathbf{A}^{-1}\|. \quad (7.38)$$

Here $\|\mathbf{A}^{-1}\|$ indicates the absolute value of the determinant of \mathbf{A}^{-1} .

Proof. Left as an exercise.

Theorem 7.18 (The Cramér-Wold Theorem) The distribution of a random vector \mathbf{y} ($p \times 1$), is completely determined by the one-dimensional distributions of all possible linear combinations of the form $\mathbf{t}'\mathbf{y}$, in which \mathbf{t} is a nonstochastic vector. The result does not assume Gaussian distributions!

Proof. Since a distribution is completely determined by its characteristic function, the Cramér-Wold Theorem can be stated $\{\phi_{\mathbf{t}'\mathbf{y}}(s) : \mathbf{t} \in \mathbb{R}^p\} \Rightarrow \phi_{\mathbf{y}}(\mathbf{t})$. \square

The characteristic function of the univariate random variable $\mathbf{t}'\mathbf{y}$ is $\phi_{\mathbf{t}'\mathbf{y}}(s) = E[\exp(i\mathbf{t}'\mathbf{y})] \forall s \in \mathbb{R}$, and the characteristic function of random vector \mathbf{y} is $\phi_{\mathbf{y}}(\mathbf{t}) = E[\exp(i\mathbf{t}'\mathbf{y})] \forall \mathbf{t} \in \mathbb{R}^p$. If $\phi_{\mathbf{y}}(\mathbf{t})$ were unknown, could we determine it from the known set $\{\phi_{\mathbf{t}'\mathbf{y}}(s) : \mathbf{t} \in \mathbb{R}^p\}$? Yes, evaluating $\phi_{\mathbf{t}'\mathbf{y}}(s)$ at $s = 1$ gives $\phi_{\mathbf{t}'\mathbf{y}}(1) = E[\exp(i\mathbf{t}'\mathbf{y})] = \phi_{\mathbf{y}}(\mathbf{t})$.

Lemma 7.3 If \mathbf{P} is a permutation matrix, then $\mathbf{P}^{-1} = \mathbf{P}'$. For a given $\mathbf{y} \in \mathbb{R}^n$, $\mathbf{z} = \mathbf{P}\mathbf{y}$ implies $\mathbf{y} = \mathbf{P}'\mathbf{z}$. Also the Jacobian of the transformation is 1.

Theorem 7.19 If \mathbf{y} is distributed with CDF $F_{\mathbf{y}}(\mathbf{y}_*)$, then permutation $\mathbf{z} = \mathbf{P}\mathbf{y}$ is distributed with CDF $F_{\mathbf{z}}(\mathbf{z}_*) = F_{\mathbf{y}}(\mathbf{P}\mathbf{z}_*)$.

Proof. Here

$$\begin{aligned} F_{\mathbf{y}}(\mathbf{P}\mathbf{z}_0) &= \Pr\{\mathbf{y} \leq \mathbf{P}\mathbf{z}_0\} = \Pr\{\mathbf{y}_* \leq \mathbf{y}_0\} \\ &= \Pr\{y_{*1} \leq y_{01}, \dots, y_{*n} \leq y_{0n}\} \\ &= \Pr\{z_{*1} \leq z_{01}, \dots, z_{*n} \leq z_{0n}\} \\ &= F_{\mathbf{z}}(\mathbf{z}_0). \end{aligned} \quad (7.39) \quad \square$$

Definition 7.11 The *Box-Cox power transformation* is given by

$$y^{(\lambda)} = \begin{cases} (y^\lambda - 1)/\lambda & \lambda \neq 0 \\ \log(y) & \lambda = 0. \end{cases} \quad (7.40)$$

The “standardized” version of $y^{(\lambda)}$ is

$$z_i^{(\lambda)} = \frac{y_i^{(\lambda)}}{\left(\prod_{i=1}^n y_i\right)^{1/n}}. \quad (7.41)$$

Box and Cox (1964, 1984) proposed maximum likelihood estimates of λ and β based on the model $y_i^{(\lambda)} \sim \mathcal{N}(\mathbf{X}_i\beta, \sigma^2)$. Others have studied the properties of the procedure (Carroll and Rupert, 1981; Hinkley and Runger, 1984).

Definition 7.12 The *Bickel-Doksum transformation* ($\lambda > 0$) is given by

$$y^{(\lambda)} = \frac{\text{sign}(y)|y|^\lambda - 1}{\lambda}. \quad (7.42)$$

The Bickel-Doksum variation of the Box-Cox transformation can cope with negative y values.

7.11 MARGINAL DISTRIBUTIONS

Definition 7.13 If an ordered set of n random variables has a joint distribution, then any subset of m of them has a joint distribution which is known by any one of the following names: the *marginal*, the *joint marginal*, the *marginal distribution*, or the *joint marginal distribution*.

Permutation matrix properties allow assuming the m variables of interest are the first m variables in the ordered set, without loss of generality. The marginal

distribution of the m variables is uniquely (a.e.) characterized by either its CDF or its characteristic function.

Theorem 7.20 If \mathbf{y} ($n \times 1$) has a distribution with CDF $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ and is partitioned as

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \quad \begin{matrix} m \times 1 \\ (n - m) \times 1 \end{matrix} \quad (7.43)$$

then the marginal joint distribution of the elements in \mathbf{y}_1 is the distribution with CDF

$$F_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) = \lim_{\mathbf{y}_{*2} \rightarrow \infty} F_{\mathbf{y}}(\mathbf{y}_{*1}, \mathbf{y}_{*2}; \boldsymbol{\theta}). \quad (7.44)$$

Proof. $F_{\mathbf{y}}(\mathbf{y}_{*1}, \mathbf{y}_{*2} = \infty; \boldsymbol{\theta}) = \Pr\{\bigcap_{j=1}^m (y_j \leq y_{j0}), \bigcap_{j=m+1}^n (y_j < \infty) | \boldsymbol{\theta}\} = F_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}). \quad \square$

Definition 7.14 For a continuous marginal distribution with CDF $F_1(\mathbf{y}_{*1}; \boldsymbol{\theta})$, the partial derivative

$$f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) = \frac{\partial^m}{\partial y_{*1} \partial y_{*2} \cdots \partial y_{*m}} F_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) \quad (7.45)$$

is the PDF of the marginal distribution if it exists $\forall \mathbf{y}_{*1} \in S$.

Theorem 7.21 If $f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta})$ exists (as just defined), then it satisfies, for $S_{\mathbf{y}_0} = \{\mathbf{y}_{*1} : \mathbf{y}_{*1} \leq \mathbf{y}_0\}$,

$$f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) = \int_{\mathbb{R}^{n-m}} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_{*2} \quad (7.46)$$

$$\int_{\mathbb{R}^m} f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) d\mathbf{y}_{*1} = 1 \quad (7.47)$$

$$f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) \geq 0 \quad \forall \mathbf{y}_{*1} \in \mathbb{R}^m \quad (7.48)$$

$$F_{\mathbf{y}_1}(\mathbf{y}_0) = \int_{S_{\mathbf{y}_0}} f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) d\mathbf{y}_{*1}. \quad (7.49)$$

Proof. Left as an exercise.

Theorem 7.22 The marginal characteristic function of \mathbf{y}_1 may be obtained by evaluating the characteristic function of $\mathbf{y}' = [\mathbf{y}'_1 \quad \mathbf{y}'_2]$ at $\mathbf{t}_2 = \mathbf{0}$,

$$\phi_{\mathbf{y}_1}(\mathbf{t}_1) = \phi_{\mathbf{y}}(\mathbf{t}) \Big|_{\mathbf{t}_2 = \mathbf{0}}. \quad (7.50)$$

Proof. $\phi_{\mathbf{y}_1}(\mathbf{t}_1) = E[\exp(i\mathbf{t}'_1 \mathbf{y}_1)] = E[\exp(i\mathbf{t}'_1 \mathbf{y}_1 + i\mathbf{0}' \mathbf{y}_2)] \equiv \phi_{\mathbf{y}}([\mathbf{t}'_1 \quad \mathbf{0}'])'. \quad \square$

7.12 INDEPENDENCE OF RANDOM VECTORS

Definition 7.15 Any pair of random vectors \mathbf{y}_1 ($n_1 \times 1$) and \mathbf{y}_2 ($n_2 \times 1$) are *statistically independent* if and only if $\Pr\{\mathbf{y}_1 \in S_1, \mathbf{y}_2 \in S_2\} = \Pr\{\mathbf{y}_1 \in S_1\}\Pr\{\mathbf{y}_2 \in S_2\} \forall$ Borel-measurable sets S_1 and S_2 .

A Borel set is “a measurable set that can be obtained from closed sets and open sets on the real line by applying the operations of union and intersection repeatedly to countable collections of sets.” (Daintith and Nelson, 1989).

Definition 7.16 The members of a set of n random vectors ($n < \infty$) $\{\mathbf{y}_j$ ($n_j \times 1$), $j \in \{1, 2, \dots, n\}\}$ are *pairwise independent* if and only if any pair of random vectors in the set are *statistically independent*.

The last two definitions do *not* imply the elements of \mathbf{y}_1 (or the elements within any other \mathbf{y}_j) are independent. An example illustrates the point. The $\{\mathbf{Y}'_i\}$ for $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ are pairwise independent. However, within a particular \mathbf{Y}'_i the elements are not independent because they have nondiagonal covariance Σ .

Definition 7.17 The members of a set of n random vectors ($n < \infty$) $\{\mathbf{y}_j$, $n_j \times 1$, $j \in \{1, 2, \dots, n\}\}$ are *mutually independent* if and only if

$$\Pr\left\{\bigcap_{j=1}^n (\mathbf{y}_j \in S_j)\right\} = \prod_{j=1}^n \Pr\{\mathbf{y}_j \in S_j\} \tag{7.51}$$

for all Borel-measurable sets $\{S_j\}$. The vectors may also be described as having as *total independence* or just *independence*.

Theorem 7.23 If $n \times 1 \mathbf{y}' = [\mathbf{y}'_1 \mathbf{y}'_2]$, with \mathbf{y}_1 $m \times 1$ and \mathbf{y}_2 $(n - m) \times 1$, has CDF $F_{\mathbf{y}}([\mathbf{y}'_{*1} \mathbf{y}'_{*2}]'; \boldsymbol{\theta})$ and characteristic function $\phi_{\mathbf{y}}([\mathbf{t}'_1 \mathbf{t}'_2]'; \boldsymbol{\theta})$ defined $\forall \mathbf{y} \in \mathbb{R}^n$, $\mathbf{t} \in \mathbb{R}^n$, then \mathbf{y}_1 and \mathbf{y}_2 are *statistically independent* \Leftrightarrow

$$F_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{y}_{*1} \\ \mathbf{y}_{*2} \end{bmatrix}; \boldsymbol{\theta}\right) = F_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{y}_{*1} \\ \infty \end{bmatrix}; \boldsymbol{\theta}\right)F_{\mathbf{y}}\left(\begin{bmatrix} \infty \\ \mathbf{y}_{*2} \end{bmatrix}; \boldsymbol{\theta}\right) \tag{7.52}$$

$\forall \mathbf{y} \in \mathbb{R}^n$. Independence requires the CDF of the joint distribution of \mathbf{y} to equal the product of the CDFs of the marginal distributions of \mathbf{y}_1 and \mathbf{y}_2 . Furthermore, \mathbf{y}_1 and \mathbf{y}_2 are *statistically independent* \Leftrightarrow

$$\phi_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \end{bmatrix}; \boldsymbol{\theta}\right) = \phi_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{t}_1 \\ \mathbf{0} \end{bmatrix}; \boldsymbol{\theta}\right)\phi_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{0} \\ \mathbf{t}_2 \end{bmatrix}; \boldsymbol{\theta}\right) \tag{7.53}$$

$\forall \mathbf{t} \in \mathbb{R}^n$. Independence requires the characteristic function of the joint distribution of \mathbf{y} to equal the product of the characteristic functions of the marginal distributions of \mathbf{y}_1 and \mathbf{y}_2 .

Proof. Left as an exercise (Cramér, 1946, p. 266, may be consulted).

Corollary 7.23.1 Similar statements of *mutual independence* can be formulated for a finite number of n subvectors $\{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_n\}$.

Proof. Left as an exercise.

Corollary 7.23.2 If the characteristic functions are replaced by MGFs, then the second part does *not* hold true for *some* distributions, even though the theorem (as stated in terms of characteristic functions) holds true for *all* distributions.

Proof. Left as an exercise. (*Hint:* Find a distribution.)

Corollary 7.23.3 Any pair of random vectors \mathbf{y}_1 ($n_1 \times 1$), \mathbf{y}_2 ($n_2 \times 1$) are *independent* if and only if $\forall \mathbf{t}_1 \in \mathbb{R}^{n_1}$ and $\forall \mathbf{t}_2 \in \mathbb{R}^{n_2}$ it follows that $\mathbf{t}'_1 \mathbf{y}_1$ is independent of $\mathbf{t}'_2 \mathbf{y}_2$. Symbolically, $\mathbf{y}_1 \perp\!\!\!\perp \mathbf{y}_2 \Leftrightarrow \mathbf{t}'_1 \mathbf{y}_1 \perp\!\!\!\perp \mathbf{t}'_2 \mathbf{y}_2 \quad \forall \mathbf{t}_1, \forall \mathbf{t}_2$.

Proof. Left as an exercise.

Corollary 7.23.4 The members of a set of n random vectors ($n < \infty$) $\{\mathbf{y}_j$ ($n_j \times 1$), $j \in \{1, 2, \dots, n\}\}$ are *pairwise independent* if and only if $\mathbf{t}'_j \mathbf{y}_j$ is independent of $\mathbf{t}'_k \mathbf{y}_k \quad \forall \mathbf{t}_j \in \mathbb{R}^{n_j}, \mathbf{t}_k \in \mathbb{R}^{n_k}$. Symbolically, $\{\mathbf{y}_j, j \in \{1, 2, \dots, n\}\}$ pairwise $\perp\!\!\!\perp \Leftrightarrow \{\mathbf{t}'_j \mathbf{y}_j \perp\!\!\!\perp \mathbf{t}'_k \mathbf{y}_k \quad \forall j \neq k, \mathbf{t}_j, \mathbf{t}_k\}$.

Corollary 7.23.5 The members of a set of n random vectors ($n < \infty$) $\{\mathbf{y}_j, n_j \times 1, j \in \{1, 2, \dots, n\}\}$ are *mutually independent* or just *independent* if and only if the members of a set of n random variables, $\{\mathbf{t}'_j \mathbf{y}_j : j \in \{1, 2, \dots, n\}\}$ are pairwise independent for all choices of $\mathbf{t}_j \in \mathbb{R}^{n_j}$.

Theorem 7.24 If $n \times 1 \mathbf{y}' = [\mathbf{y}'_1 \mathbf{y}'_2]$, with \mathbf{y}_1 $m \times 1$ and \mathbf{y}_2 $(n - m) \times 1$, has PDF $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ and the marginal distribution of \mathbf{y}_j has PDF $f_j(\mathbf{y}_{*j}; \boldsymbol{\theta}), j \in \{1, 2\}$, then \mathbf{y}_1 and \mathbf{y}_2 are *independent* if and only if $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) = f_1(\mathbf{y}_{*1}; \boldsymbol{\theta})f_2(\mathbf{y}_{*2}; \boldsymbol{\theta}) \quad \forall \mathbf{y} \in \mathbb{R}^n$.

Proof. Left as an exercise.

7.13 CONDITIONAL DISTRIBUTIONS

Definition 7.18 If $n \times 1 \mathbf{y}' = [\mathbf{y}'_1 \mathbf{y}'_2]$, with \mathbf{y}_1 $m \times 1$ and \mathbf{y}_2 $(n - m) \times 1$, has a joint distribution with CDF $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$, the marginal of \mathbf{y}_1 has CDF $F_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta})$, $S \subset \mathbb{R}^m$ is a Borel-measurable set, and a function $F_{2|1}(\mathbf{y}_{*2} | \mathbf{y}_{01}; \boldsymbol{\theta}_c)$ exists such that

$$\Pr\{\mathbf{y}_1 \in S, \mathbf{y}_2 \leq \mathbf{y}_{02}\} = \int_S F_{2|1}(\mathbf{y}_{*2}|\mathbf{y}_{*1}; \boldsymbol{\theta}_c) dF_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta}) \quad (7.54)$$

$\forall S \subset \mathbb{R}^p$ and $\mathbf{y}_{02} \in \mathbb{R}^{n-1}$, then $F_{2|1}(\mathbf{y}_{*2}|\mathbf{y}_{*1}; \boldsymbol{\theta}_c)$ is called the *CDF of the conditional distribution* of the random vector $\mathbf{y}_2|\mathbf{y}_1 = \mathbf{y}_{01}$. Furthermore any random vector with CDF $F_{2|1}(\mathbf{y}_{*2}|\mathbf{y}_{01}; \boldsymbol{\theta}_c)$ may be called $\mathbf{y}_2|\mathbf{y}_1 = \mathbf{y}_{01}$.

In practice, $\mathbf{y}_2|\mathbf{y}_1 = \mathbf{y}_{01}$ and $\mathbf{y}_2|\mathbf{y}_1$ are often used interchangeably, although they do not represent precisely the same random vector.

Theorem 7.25 If the joint distribution of $n \times 1 \mathbf{y}' = [\mathbf{y}'_1 \mathbf{y}'_2]$, with \mathbf{y}_1 $m \times 1$ and \mathbf{y}_2 $(n - m) \times 1$, is absolutely continuous with PDF $f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ and the marginal distribution of \mathbf{y}_1 has PDF $f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta})$ with support $S_1 \subseteq \mathbb{R}^m$ (i.e., $\Pr\{\mathbf{y}_1 \in S_1\} = 1$), then

(a) the CDF of the conditional distribution of $\mathbf{y}_2|\mathbf{y}_1 = \mathbf{y}_{01}$ is, with $S_{\mathbf{y}_{02}} = \{\mathbf{y}_{*2} : \mathbf{y}_{*2} \leq \mathbf{y}_{02}\}$,

$$F_{2|1}(\mathbf{y}_{02}|\mathbf{y}_1; \boldsymbol{\theta}_c) = \int_{S_{\mathbf{y}_{02}}} \frac{f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})}{f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta})} d\mathbf{y}_{*2} \quad (7.55)$$

and (b) the density of the conditional distribution is, $\forall \mathbf{y}_{*2} \in S_1$,

$$f_{2|1}(\mathbf{y}_{*2}|\mathbf{y}_{*1}; \boldsymbol{\theta}_c) = \frac{f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})}{f_{\mathbf{y}_1}(\mathbf{y}_{*1}; \boldsymbol{\theta})}. \quad (7.56)$$

In many cases conditional distributions do not exist. A conditional probability exists only if the probability in the denominator is greater than zero.

7.14 (JOINT) MOMENTS OF MULTIVARIATE DISTRIBUTIONS

Definition 7.19 (a) The *expected value* of random vector \mathbf{y} is a vector-valued function of the joint CDF $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$ defined as

$$\mathbf{E}(\mathbf{y}) = \begin{bmatrix} \mathbf{E}(y_1) \\ \vdots \\ \mathbf{E}(y_n) \end{bmatrix} \quad (7.57)$$

with

$$\mathbf{E}(y_j) = \begin{cases} \int_{\mathbb{R}^n} y_{*j} dF_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) & \text{(in general)} \\ \int_{\mathbb{R}^n} y_{*j} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_* & \text{(if the PDF exists).} \end{cases} \quad (7.58)$$

(b) The value of $E(\mathbf{y})$ is an $n \times 1$ vector $\boldsymbol{\mu}$ also known as the *mean vector*. The vector of means exists if all its elements $\{\mu_j\}$ are finite.

Definition 7.20 If \mathbf{y} is a random vector which follows a distribution with CDF $F_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta})$ and

$$\mathbf{h}(\mathbf{y}) = \begin{bmatrix} h_1(\mathbf{y}) \\ \vdots \\ h_k(\mathbf{y}) \end{bmatrix} \tag{7.59}$$

is a vector of functions of $\mathbf{y} \in \mathbb{R}^n$ which are integrable with respect to $F_{\mathbf{y}}(\mathbf{y}; \boldsymbol{\theta})$, then the *expected value* of $E[\mathbf{h}(\mathbf{y})]$ is defined as the following elementwise integral operation:

$$E[\mathbf{h}(\mathbf{y})] = \begin{bmatrix} E[h_1(\mathbf{y})] \\ \vdots \\ E[h_k(\mathbf{y})] \end{bmatrix} \tag{7.60}$$

with

$$E[h_j(\mathbf{y})] = \begin{cases} \int_{\mathbb{R}^n} h_j(\mathbf{y}_*) dF_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) & \text{(in general)} \\ \int_{\mathbb{R}^n} h_j(\mathbf{y}_*) f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_* & \text{(if the PDF exists).} \end{cases} \tag{7.61}$$

The moments exist if the integrals are finite.

Definition 7.21 (a) If \mathbf{y} follows a distribution with CDF $F_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta})$, $\mathbf{c} \in \mathbb{R}^k$, $m > 0$, and

$$\mathbf{h}(\mathbf{y}) = \begin{bmatrix} h_1(\mathbf{y}) \\ \vdots \\ h_k(\mathbf{y}) \end{bmatrix} = \begin{bmatrix} (y_1 - c_1)^m \\ \vdots \\ (y_k - c_k)^m \end{bmatrix}, \tag{7.62}$$

then $E[\mathbf{h}(\mathbf{y})] = [E[h_1(\mathbf{y})] \cdots E[h_k(\mathbf{y})]]'$ with

$$E[h_j(\mathbf{y})] = \begin{cases} \int_{\mathbb{R}^n} (y_{*j} - c_j)^m dF_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) & \text{(in general)} \\ \int_{\mathbb{R}^n} (y_{*j} - c_j)^m f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_* & \text{(if the PDF exists).} \end{cases} \tag{7.63}$$

(b) If $\mathbf{c} = \mathbf{0}$, then $E[\mathbf{h}(\mathbf{y})] \equiv [\mu_{m',1} \cdots \mu_{m',k}]' \equiv \boldsymbol{\mu}_{m'}$ is the set of *moments about $\mathbf{0}$* of order m ,

(c) If $\mathbf{c} = \boldsymbol{\mu}_{1'}$, then $E[\mathbf{h}(\mathbf{y})] \equiv [\mu_{m,1} \cdots \mu_{m,k}]' \equiv \boldsymbol{\mu}_m$ is the set of *central moments* of order m .

Usually we work with random *vectors* because distributions, covariance matrices, and many related properties of scalar random variables remain well-

defined for random vectors. In addition, we sometimes consider random matrices, such as an $n \times p$ data matrix \mathbf{Y} . Similarly, an estimated covariance (dispersion) matrix is a symmetric random matrix.

Definition 7.22 If $\mathbf{y} = \{y_{jk}\}$ is an $n \times p$ array of random variables having a joint distribution and $E(y_{jk})$ exists (finitely) for all (j, k) , the $n \times p$ matrix of expected values is $E(\mathbf{Y}) = \{E(y_{jk})\}$, with expectation an elementwise operation.

Theorem 7.26 If \mathbf{Y} and \mathbf{X} are random $n \times p$ arrays and \mathbf{A} , \mathbf{B} , and \mathbf{C} are constant matrices which conform for the matrix operations implied, then

$$E(\mathbf{X} + \mathbf{Y}) = E(\mathbf{X}) + E(\mathbf{Y}) \quad (7.64)$$

$$E(\mathbf{X} + \mathbf{C}) = E(\mathbf{X}) + \mathbf{C} \quad (7.65)$$

$$E(\mathbf{C}) = \mathbf{C} \quad (7.66)$$

$$E(\mathbf{A}\mathbf{X}) = \mathbf{A}E(\mathbf{X}) \quad (7.67)$$

$$E(\mathbf{X}\mathbf{B}) = [E(\mathbf{X})]\mathbf{B}. \quad (7.68)$$

Definition 7.23 If it exists, the *covariance* or *dispersion* matrix $\mathcal{V}(\mathbf{y})$ is

$$\mathcal{V}(\mathbf{y}) = E[(\mathbf{y} - \boldsymbol{\mu})(\mathbf{y} - \boldsymbol{\mu})'], \quad (7.69)$$

with $\boldsymbol{\mu} = E(\mathbf{y})$ and the (j, k) element a *variance* $\mathcal{V}(y_j)$ if $j = k$ and a *covariance* $\mathcal{V}(y_j, y_k)$ otherwise. In particular,

$$E[(y_j - \mu_j)(y_k - \mu_k)] = \begin{cases} \int (y_{*j} - \mu_j)(y_{*k} - \mu_k) dF_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) & \text{(in general)} \\ \int (y_{*j} - \mu_j)(y_{*k} - \mu_k) f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\theta}) d\mathbf{y}_* & \text{(if a.c.).} \end{cases} \quad (7.70)$$

The array of second-order moments exists if and only if all its elements $\{\sigma_{jk}\}$ have finite values. Although using σ_j^2 in place of σ_{jj} seems tempting, consistent use of σ_{jj} helps prevent confusion.

Theorem 7.27 If random vector $\mathbf{y} = [y_1 \cdots y_n]'$ has a continuous joint distribution, $E(\mathbf{y}) = \boldsymbol{\mu}$, $\mathcal{V}(\mathbf{y}) = \boldsymbol{\Sigma}$, and $\sigma_{jj} = \mathcal{V}(y_j) < \infty \forall j \in \{1, 2, \dots, n\}$, then, for conforming constants \mathbf{A} , \mathbf{a} , and \mathbf{t} , the following all hold.

(a) $(\sigma_{jk})^2 = [\mathcal{V}(y_j, y_k)]^2 \leq \mathcal{V}(y_j)\mathcal{V}(y_k) < \infty$

(b) $\mathcal{V}(\mathbf{y} + \mathbf{a}) = \mathcal{V}(\mathbf{y}) = \boldsymbol{\Sigma} \forall \mathbf{a} \in \mathfrak{R}^n$

(c) $\mathcal{V}(\mathbf{A}\mathbf{y}) = \mathbf{A}\mathcal{V}(\mathbf{y})\mathbf{A}' = \mathbf{A}\boldsymbol{\Sigma}\mathbf{A}' \forall \mathbf{A} \in \mathfrak{R}^{m \times n}$

(d) $\mathcal{V}(\mathbf{t}'\mathbf{y}) = \mathbf{t}'\mathcal{V}(\mathbf{y})\mathbf{t} = \mathbf{t}'\boldsymbol{\Sigma}\mathbf{t} \geq 0 \forall \mathbf{t} \in \mathfrak{R}^n$

(e) $\boldsymbol{\Sigma}$ is either positive definite or positive semidefinite

(f) $\boldsymbol{\Sigma} = E[(\mathbf{y} - \mathbf{d})(\mathbf{y} - \mathbf{d})'] - (\boldsymbol{\mu} - \mathbf{d})(\boldsymbol{\mu} - \mathbf{d})' \forall \mathbf{d} \in \mathfrak{R}^n$

(g) $\boldsymbol{\Sigma} = E(\mathbf{y}\mathbf{y}') - \boldsymbol{\mu}\boldsymbol{\mu}'$ for $\mathbf{d} = \mathbf{0}$

Proof. Left as an exercise. *Hints:* The Cauchy-Schwartz inequality, $[\int_a^b f(x)g(x)dx]^2 \leq \int_a^b [f(x)]^2 dx \int_a^b [g(x)]^2 dx, \Rightarrow$ (a). The integral definition helps with (b). Also (c) \Rightarrow (d) \Leftrightarrow (e).

Theorem 7.28 If $\mathbf{y}' = [\mathbf{y}'_1 \mathbf{y}'_2]$, with \mathbf{y}_j of dimension p_j , and all elements have finite second moments, then

$$\mathcal{V}(\mathbf{y}) = \Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \tag{7.71}$$

with Σ_{jk} ($p_j \times p_k$) described as a covariance matrix. Also $\Sigma_{jj} = \mathcal{V}(\mathbf{y}_j)$. Here $\Sigma_{12} = \Sigma'_{21} = E[(\mathbf{y}_1 - \boldsymbol{\mu}_1)(\mathbf{y}_2 - \boldsymbol{\mu}_2)']$, often written $\mathcal{V}(\mathbf{y}_1, \mathbf{y}_2) = [\mathcal{V}(\mathbf{y}_2, \mathbf{y}_1)]'$. For (conformable) constants \mathbf{A} and \mathbf{B} , $\mathcal{V}()$ is a bilinear operator,

$$\mathcal{V}(\mathbf{A}\mathbf{y}_1, \mathbf{B}\mathbf{y}_2) = \mathbf{A}\mathcal{V}(\mathbf{y}_1, \mathbf{y}_2)\mathbf{B}', \tag{7.72}$$

and $\mathcal{V}()$ is invariant to adding constants,

$$\mathcal{V}(\mathbf{a} + \mathbf{A}\mathbf{y}_1, \mathbf{b} + \mathbf{B}\mathbf{y}_2) = \mathbf{A}\mathcal{V}(\mathbf{y}_1, \mathbf{y}_2)\mathbf{B}'. \tag{7.73}$$

Proof. Left as an exercise.

Theorem 7.29 For a random vector \mathbf{z} following any distribution with finite second moments and \mathbf{T} a fixed and conforming matrix, $E(\mathbf{T}\mathbf{z}) = \mathbf{T}\boldsymbol{\mu}_z$ and $\mathcal{V}(\mathbf{T}\mathbf{z}) = \mathbf{T}\Sigma_z\mathbf{T}'$.

Definition 7.24 If $n \times 1$ random vector \mathbf{y} has $\mathcal{V}(\mathbf{y}) = \Sigma = \{\sigma_{jk}\}$ with $0 < \sigma_{jj} < \infty$, the $n \times n$ correlation matrix is

$$\mathbf{P} = [\text{Dg}(\{\sigma_{11}, \dots, \sigma_{nn}\})]^{-1/2} \Sigma [\text{Dg}(\{\sigma_{11}, \dots, \sigma_{nn}\})]^{-1/2}, \tag{7.74}$$

with (j, k) element $\rho_{jk} = \sigma_{jk}(\sigma_{jj}\sigma_{kk})^{-1/2}$. If $j = k$, then $\rho_{jj} = 1$, and otherwise ρ_{jk} is the correlation between y_j and y_k .

Theorem 7.30 If $n \times 1$ random vector \mathbf{y} has $\mathcal{V}(\mathbf{y}) = \Sigma$, with $0 < \sigma_{jj} < \infty$, then

- (a) $\rho_{jk} \in [-1, 1] \forall j, k$
- (b) \mathbf{P} ($n \times n$) is symmetric and nonnegative definite.
- (c) For any given Σ the unique value of \mathbf{P} is

$$\mathbf{P} = [\text{Dg}(\Sigma)]^{-1/2} \Sigma [\text{Dg}(\Sigma)]^{-1/2}. \tag{7.75}$$

(d) For any given \mathbf{P} the value of Σ depends on knowing $\text{Dg}(\Sigma)$ because

$$\Sigma = [\text{Dg}(\Sigma)]^{1/2} \mathbf{P} [\text{Dg}(\Sigma)]^{1/2}. \tag{7.76}$$

Proof. (a) The Cauchy-Schwartz inequality implies $\sigma_{jk}^2 \leq \sigma_{jj}\sigma_{kk}$, which implies $\rho_{jk} \in [-1, 1]$. (b) Σ is positive semidefinite and congruent to \mathbf{P} . Hence

P is also positive semidefinite (by Sylvester's law of inertia). (c) is true by definition. (d) Any valid choice of standard deviations (strictly positive and finite) yields a valid value for Σ .

Definition 7.25 For

$$\mathcal{V} \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \tag{7.77}$$

(a) the vectors \mathbf{y}_1 and \mathbf{y}_2 are *uncorrelated* $\Leftrightarrow \Sigma_{12} = \mathbf{0} \Leftrightarrow P_{12} = \mathbf{0}$.

(b) The elements of \mathbf{y}_j are *uncorrelated* $\Leftrightarrow \Sigma_{jj}$ is diagonal $\Leftrightarrow P_{jj} = I$.

A lack of correlation does not necessarily imply statistical independence.

7.15 CONDITIONAL MOMENTS OF DISTRIBUTIONS

Here we consider properties of $\mathbf{y}|\mathbf{x}$. In particular, $E[E(\mathbf{y}|\mathbf{x})] = E(\mathbf{y})$, while $E[\mathcal{V}(\mathbf{y}|\mathbf{x})] \leq \mathcal{V}(\mathbf{y})$. The second result should seem reasonable because $\mathcal{V}(\mathbf{y}|\mathbf{x} = \mathbf{x}_0) \leq \mathcal{V}(\mathbf{y}) \forall \mathbf{x}_0$, with equality whenever $\mathbf{y} \perp\!\!\!\perp \mathbf{x}$. Knowing the specific realization of \mathbf{x} can only reduce uncertainty about the value of \mathbf{y} .

Lemma 7.4 If \mathbf{y} and \mathbf{x} have a joint distribution, then the CDF and PDF (if it exists) are

$$\begin{aligned} F_{\mathbf{y}}(\mathbf{y}_*) &= E[F_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_*|\mathbf{x}_*)] \quad (\text{always}) \\ f_{\mathbf{y}}(\mathbf{y}_*) &= E[f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_*|\mathbf{x}_*)] \quad (\text{if it exists}). \end{aligned} \tag{7.78}$$

Proof. Left as an exercise.

Theorem 7.31 If \mathbf{y} and \mathbf{x} have a joint distribution then $E(\mathbf{y}) = E[E(\mathbf{y}|\mathbf{x})]$, which for clarity may be written $E(\mathbf{y}) = E_{\mathbf{x}}[E_{\mathbf{y}}(\mathbf{y}|\mathbf{x})]$.

Proof. (For absolutely continuous distributions.)

$$\begin{aligned} E(\mathbf{y}) &= \int \mathbf{y}_* f_{\mathbf{y}}(\mathbf{y}_*) d\mathbf{y}_* \\ &= \int \mathbf{y}_* \left[\int f_{\mathbf{y},\mathbf{x}}(\mathbf{y}_*, \mathbf{x}_*) d\mathbf{x}_* \right] d\mathbf{y}_* \\ &= \int \int \mathbf{y}_* f_{\mathbf{y},\mathbf{x}}(\mathbf{y}_*|\mathbf{x}_*) f_{\mathbf{x}}(\mathbf{x}_*) d\mathbf{x}_* d\mathbf{y}_* \\ &= \int \left[\int \mathbf{y}_* f_{\mathbf{y}|\mathbf{x}}(\mathbf{y}_*|\mathbf{x}_*) d\mathbf{y}_* \right] f_{\mathbf{x}}(\mathbf{x}_*) d\mathbf{x}_* \\ &= E[E(\mathbf{y}|\mathbf{x})] = E_{\mathbf{x}}[E_{\mathbf{y}}(\mathbf{y}|\mathbf{x})]. \end{aligned} \tag{7.79}$$

□

Theorem 7.32 If \mathbf{y} and \mathbf{x} have a joint distribution and the moments exist, then

$$\mathcal{V}(\mathbf{y}) = E[\mathcal{V}(\mathbf{y}|\mathbf{x})] + \mathcal{V}[E(\mathbf{y}|\mathbf{x})]. \quad (7.80)$$

Proof. If $\boldsymbol{\mu} = E(\mathbf{y})$ and $\boldsymbol{\mu}_c = E(\mathbf{y}|\mathbf{x})$ then

$$\begin{aligned} \mathcal{V}(\mathbf{y}) &= E(\mathbf{y}\mathbf{y}') - \boldsymbol{\mu}\boldsymbol{\mu}' \\ &= [E(\mathbf{y}\mathbf{y}') - E(\boldsymbol{\mu}_c\boldsymbol{\mu}_c')] + [E(\boldsymbol{\mu}_c\boldsymbol{\mu}_c') - \boldsymbol{\mu}\boldsymbol{\mu}']. \end{aligned} \quad (7.81)$$

Also

$$\begin{aligned} E(\mathbf{y}\mathbf{y}') - E(\boldsymbol{\mu}_c\boldsymbol{\mu}_c') &= E\{E(\mathbf{y}\mathbf{y}'|\mathbf{x}) - E(\mathbf{y}|\mathbf{x})[E(\mathbf{y}|\mathbf{x})]'\} \\ &= E[\mathcal{V}(\mathbf{y}|\mathbf{x})] \end{aligned} \quad (7.82)$$

and

$$\begin{aligned} E(\boldsymbol{\mu}_c\boldsymbol{\mu}_c') - \boldsymbol{\mu}\boldsymbol{\mu}' &= E(\boldsymbol{\mu}_c\boldsymbol{\mu}_c') - E(\boldsymbol{\mu}_c)E(\boldsymbol{\mu}_c') \\ &= \mathcal{V}(\boldsymbol{\mu}_c) \\ &= \mathcal{V}[E(\mathbf{y}|\mathbf{x})]. \end{aligned} \quad (7.83) \quad \square$$

7.16 SPECIAL CONSIDERATIONS FOR RANDOM MATRICES

At first glance, generalizing to matrices introduces no great complication. It obviously makes sense to define, for $\mathbf{X} = \{x_{jk}\}$ with x_{jk} random,

$$E(\mathbf{X}) = \{E(x_{jk})\}. \quad (7.84)$$

However, even when $\{x_{jk}\}$ have finite second moments, $\mathcal{V}(\mathbf{X})$ has no obvious meaning, although $\mathcal{V}[\text{vec}(\mathbf{X})]$ and $\{\mathcal{V}(x_{jk})\}$ are well defined (and typically quite different). For probability calculations and other operations directly involving matrix elements, expressions in terms of the $\text{vec}()$ and $\text{vech}()$ operators seem most useful. In contrast, it is often advantageous to consider moment calculations in terms of the original matrix form. Consequently the characteristic (and moment and cumulant generating) functions have been generalized.

The characteristic function tends to be used more in the study of random vectors than for scalars. In turn, it plays an even bigger role for random matrices. Scalar (continuous) random variables that are commonly studied rarely fail to have a density. In contrast, singular vector and matrix Gaussian random arrays, which do not have a density, arise naturally in the study of data analysis. Even when a density for a vector or matrix exists, it may not have a convenient closed form. Useful expressions for CDFs are even more rare. Furthermore, even when known, densities for random vectors and matrices often prove difficult to manipulate, for both analytic and computational purposes.

At least for the random vectors and matrices of interest in the present book, the characteristic functions are never more complicated than corresponding densities (when they exist), and often much simpler. In turn, much simpler proofs can often

be found in terms of characteristic functions. Most importantly, the proofs often apply to a wider range of distributions by tolerating random vectors and matrices which do not have densities (due to singularities in the distribution function).

The first task is to generalize the definition of the characteristic function of a vector to a matrix. The definition follows the one give by Gupta and Nagar (2000, p45–46, with some notation changed).

Definition 7.26 With $i = \sqrt{-1}$, T an arbitrary real $n \times p$ matrix, and random $\{x_{jk}\}$ having a well-defined joint distribution, the *characteristic function* of $n \times p$ $\mathbf{X} = \{x_{jk}\}$ is

$$\phi_{\mathbf{X}}(T) = E\{\exp[\text{tr}(iT'X)]\}. \quad (7.85)$$

Definition 7.27 With T an arbitrary real $n \times p$ matrix and random $\{x_{jk}\}$ having a well-defined joint distribution, the *moment generating function* of $n \times p$ $\mathbf{X} = \{x_{jk}\}$, when it exists, is

$$m_{\mathbf{X}}(T) = E\{\exp[\text{tr}(T'X)]\}. \quad (7.86)$$

Definition 7.28 With T an arbitrary real $n \times p$ matrix, and random $\{x_{jk}\}$ having a well-defined joint distribution, the *cumulant generating function* of $n \times p$ $\mathbf{X} = \{x_{jk}\}$, when it exists, is

$$c_{\mathbf{X}}(T) = \log[m_{\mathbf{X}}(T)]. \quad (7.87)$$

The following lemma generalizes the result about a linear transformation of a random vector. The matrix form allows simple proofs of many results concerning Gaussian variables and corresponding covariance estimators. Most importantly, the lemma allows less than full rank transformations, even when applied to random variables lacking a density. In applications in later chapters, the resulting characteristic function is often seen to correspond to a known distribution.

Lemma 7.5 If $Y = AXB + C$ for random X and finite constants $\{A, B, C\}$ of any rank and conforming size $(n_1 \times p_1) = (n_1 \times n)(n \times p)(p \times p_1) + (n_1 \times p_1)$, then

$$\phi_Y(T) = \phi_X(A'TB') \cdot \exp[i \text{tr}(T'C)]. \quad (7.88)$$

Proof. The result is proven by writing

$$\begin{aligned} E(\exp\{\text{tr}[iT'(AXB + C)]\}) &= E\{\exp[\text{tr}(iT'AXB) + \text{tr}(iT'C)]\} \\ &= E\{\exp[\text{tr}(iBT'AX)]\exp[\text{tr}(iT'C)]\} \quad (7.89) \\ &= E(\exp\{\text{tr}[i(A'TB)'X]\})\exp[i \text{tr}(T'C)]. \quad \square \end{aligned}$$

CHAPTER 8

Scalar, Vector, and Matrix Gaussian Distributions

8.1 MOTIVATION

The Gaussian distribution, which was discovered by de Moivre in 1733, has likely been studied more than any other. The distribution often provides a plausible model in a wide range of applications, primarily because it arises so often as the asymptotic distribution specified in central limit theorems. Equally importantly, Gaussian distributions and distributions of functions of Gaussian variables, such as chi squares, typically provide the backbone of many distribution-free results for large samples. The statement holds for many scalar, vector, and matrix forms.

In the context of linear models, assuming Gaussian error terms typically implies estimators of means and treatment effects (expected-value parameters) often follow Gaussian distributions. Gaussian distributions may provide an appropriate model of a first moment (mean, location) in scalar, vector, or matrix form.

We avoid the term “normal” distribution to avoid the implication that the distribution should be expected to apply or is “usual.” The term originated from a philosophical position about universal laws of nature. Referring to a “Gaussian” distribution helps emphasize it is one particular distribution among many and that the specific choice must be justified in any particular application. We mostly ignore the important practical problem of defending the use of the Gaussian assumption for a particular set of data and proceed under the assumption.

Statisticians have devised many ways of characterizing the Gaussian distribution. The approaches vary in the simplicity of assumptions and presentation, generality of parameters allowed, and ease of learning. We maintain complete generality but sacrifice some simplicity of presentation to make learning easier. We begin with the standard scalar Gaussian and define the general scalar Gaussian as a (scalar) linear transformation of the standard Gaussian. Some transformation constants lead to nonsingular distribution functions (they have continuous derivatives which provide a density) and correspond to nondegenerate random variables. Other transformation constants lead to distribution functions which are singular (they do not have continuous derivatives or a density) and correspond to degenerate random variables.

The same sequence of constructions generalizes from scalars to vectors, beginning with a vector of independent standard Gaussian variables. In turn, many of the convenient and special properties of the vector Gaussian are derived. The chapter concludes with consideration of matrix forms of Gaussian variables. Such forms are particularly useful in the study of multivariate models, especially multivariate linear models.

8.2 THE SCALAR GAUSSIAN DISTRIBUTION

We delay defining the scalar Gaussian until we prove the function we claim is its density is in fact a density. Lemma 8.1 provides the standard Gaussian result.

Lemma 8.1

$$\int_{-\infty}^{\infty} \exp(-z_*^2/2) dz_* = (2\pi)^{1/2}. \quad (8.1)$$

Proof. The lemma is true if and only if the square of the integral is 2π . Here

$$\begin{aligned} \left[\int_{-\infty}^{\infty} \exp(-z_*^2/2) dz_* \right]^2 &= \left[\int_{-\infty}^{\infty} \exp(-x_*^2/2) dx_* \right] \left[\int_{-\infty}^{\infty} \exp(-y_*^2/2) dy_* \right] \\ &= \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp[-(x_*^2 + y_*^2)/2] dx_* dy_*. \end{aligned} \quad (8.2)$$

Switching to polar coordinates gives $x_* = r_* \cos(\theta_*)$ and $y_* = r_* \sin(\theta_*)$, with $r_*^2 = x_*^2 + y_*^2$ and $dx_* dy_* = r_* dr_* d\theta_*$ (r_* is the Jacobian of the transformation). In turn,

$$\begin{aligned} \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \exp[-(x_*^2 + y_*^2)/2] dx_* dy_* &= \int_0^{2\pi} \int_0^{\infty} \exp(-r_*^2/2) r_* dr_* d\theta_* \\ &= \left(\int_0^{2\pi} d\theta_* \right) \int_0^{\infty} \exp(-r_*^2/2) r_* dr_* \\ &= 2\pi \left[-\int_{r_*=0}^{\infty} \exp(-r_*^2/2) \right] = 2\pi. \end{aligned} \quad (8.3) \quad \square$$

Theorem 8.1 The function $f_z(z_*) = (2\pi)^{-1/2} \exp(-z_*^2/2)$ is the density of an a.c. random variable with domain the real line.

Proof. The function is nonnegative and integrates to 1 over the real line (by the lemma). By Theorem 7.1, the function is a density of a random variable. \square

Corollary 8.1.1 For $f_z(z_*) = (2\pi)^{-1/2} \exp(-z_*^2/2)$,

- (a) the MGF of z is $m_z(t) = \exp(-t^2/2)$ and
- (b) the CF of z is $\phi_z(t) = \exp(-t^2/2)$.
- (c) $E(z) = 0$ and $\mathcal{V}(z) = E(z^2) = 1$.
- (d) All odd moments are zero.
- (e) If $m \in \{2, 4, \dots\}$, then $E(z^m) = (m - 1)(m - 3) \cdots (3)(1)$.
- (f) The cumulants $\{\kappa_m\}$ are all zero for $m > 2$.

Proof. (a)

$$\begin{aligned}
 E[\exp(zt)] &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp(z_*t) \exp(-z_*^2/2) dz_* \\
 &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[(-z_*^2 + 2z_*t - t^2 + t^2)/2] dz_* \\
 &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[-(z_* - t)^2/2 + t^2/2] dz_* \\
 &= \exp(t^2/2) (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[-(z_* - t)^2/2] dz_* \\
 &= \exp(t^2/2) (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[-(u_*)^2/2] du_* \\
 &= \exp(t^2/2) \cdot 1.
 \end{aligned} \tag{8.4}$$

(b) The characteristic function of z is

$$\begin{aligned}
 E[\exp(izt)] &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp(iz_*t) \exp(-z_*^2/2) dz_* \\
 &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[(-z_*^2 + 2iz_*t + t^2 - t^2)/2] dz_* \\
 &= (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[-(z_* - it)^2/2 - t^2/2] dz_* \\
 &= \exp(-t^2/2) (2\pi)^{-1/2} \int_{-\infty}^{\infty} \exp[-(z_* - it)^2/2] dz_* \\
 &\quad \text{(skipping a complex integration)} \\
 &= \exp(-t^2/2) \cdot 1.
 \end{aligned} \tag{8.5}$$

- (c) Moments are easily computed by integration or from derivatives of the MGF.
- (d) True due to symmetry about the origin.
- (e) and (f) are left as exercises (consider induction). □

Corollary 8.1.2 If the random variable z has density $f_z(z_*) = (2\pi)^{-1/2} \exp(-z_*^2/2)$ and domain \mathfrak{R}^1 , then $y = \sigma z + \mu$,

- (a) for (real) constants $-\infty < \mu < \infty$ and $0 < \sigma^2 < \infty$, has domain \mathfrak{R}^1 and density

$$f_y(y_*; \mu, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp[-(y_* - \mu)^2/(2\sigma^2)]. \tag{8.6}$$

- (b) If $\sigma^2 = 0$ and $-\infty < \mu < \infty$, then $y = \sigma z + \mu$ is a discrete and degenerate random variable, with $\Pr\{y = \mu\} = 1$, and y does not have a density.

Proof. If $0 < \sigma^2 < \infty$, the function must be nonnegative, integrate to 1 over the domain, and integrate to provide the CDF of y for $0 < \sigma^2 < \infty$. The function $f_y(y_*; \mu, \sigma^2)$ is clearly nonnegative for all y_* (finite μ , and $0 < \sigma^2 < \infty$). Also,

$$\begin{aligned} \Pr\{y \leq y_0\} &= \Pr\{(\sigma z + \mu) \leq y_0\} \\ &= \Pr\{z \leq (y_0 - \mu)/\sigma\} \\ &= \int_{-\infty}^{(y_0 - \mu)/\sigma} (2\pi)^{-1/2} \exp(-z_*^2/2) dz_*. \end{aligned} \quad (8.7)$$

If $y_0 \rightarrow \infty$ then the integral is 1 over the real line (by the lemma). The transformation $t = \sigma z_* + \mu$ gives $dt_* = \sigma dz_*$ and $z_* = (t_* - \mu)/\sigma$. In turn,

$$\begin{aligned} \Pr\{y \leq y_0\} &= \int_{-\infty}^{y_0} (2\pi)^{-1/2} \exp[-(t_* - \mu)^2 / (2\sigma^2)] \sigma^{-1} dt_* \\ &= \int_{-\infty}^{y_0} f_y(t_*; \mu, \sigma^2) dt_*. \end{aligned} \quad (8.8)$$

If $\sigma^2 = 0$ then $y = 0 \cdot z + \mu \equiv \mu$ and $\Pr\{y = \mu\} = 1$. If so, y is a discrete random variable and does not have a density. \square

Corollary 8.1.3 For (real) constants $-\infty < \mu < \infty$ and $0 \leq \sigma^2 < \infty$, the characteristic function of $y = \sigma z + \mu$ is

$$\begin{aligned} \phi_y(t) &= e^{it\mu} \phi_z(\sigma t) \\ &= \exp(it\mu - \sigma^2 t^2 / 2), \end{aligned} \quad (8.9)$$

and the MGF is $m_y(t) = \exp(t\mu + \sigma^2 t^2 / 2)$. Furthermore, $E(y) = \mu$ and $\mathcal{V}(y) = \sigma^2$, which fully characterize the distribution.

Definition 8.1 (a) A scalar (real) random variable z with density

$$f_z(z_*) = (2\pi)^{-1/2} \exp(-z_*^2/2) \quad (8.10)$$

is said to follow a *standard Gaussian* distribution, written $z \sim \mathcal{N}(0, 1)$.

(b) If μ and $0 < \sigma^2$ are finite real constants, then the scalar random variable $y = \sigma z + \mu$ has density

$$f_y(y_*; \mu, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp[-(y_* - \mu)^2 / (2\sigma^2)] \quad (8.11)$$

and is said to follow a *Gaussian distribution*, written $y \sim \mathcal{N}(\mu, \sigma^2)$.

(c) If $\sigma^2 = 0$ (and μ is finite), then y does not have a density, and follows a *singular Gaussian distribution*, written $y \sim \mathcal{SN}(\mu, 0)$, with $\Pr\{y = \mu\} = 1$.

(d) Writing $y \sim (\mathcal{S})\mathcal{N}(\mu, \sigma^2)$ indicates y may be either singular or nonsingular (with $0 \leq \sigma^2 < \infty$), which is sometimes written $y \sim \mathcal{N}(\mu, \sigma^2)$.

The adjective “singular” reflects a singularity at $\sigma^2 = 0$ for $f_y(y; \mu, \sigma^2)$, considered as a function of σ^2 . The corresponding CDF, considered as a function

of y_* , has a singularity at μ . Although $y \sim \mathcal{SN}(\mu, 0)$ does not have a density, its CDF and other properties are well defined. The singular Gaussian can be appropriately treated as the limiting case of a (nonsingular) Gaussian for $\sigma^2 \rightarrow 0$.

Other singularities arise if $\mu \rightarrow \pm\infty$ or $\sigma^2 \rightarrow +\infty$ (separately). Although such conditions arise naturally as limiting cases, it is very common, although not universal, to exclude them from the definition of singular Gaussian, as we do here. In contrast, electrical engineers sometimes find it convenient to describe a theoretical source of random noise as having equal power at all frequencies. Such “white” noise corresponds to a Gaussian random variable with infinite σ^2 .

No other values for $\{\mu, \sigma^2\}$ lead to a valid distribution (at least without strong side conditions). In particular, if $-\infty < \sigma^2 < 0$ and $-\infty < \mu < \infty$, then the integral of $\exp[-(y_* - \mu)^2 / (2\sigma^2)]$ over the real line does not exist. Furthermore, having μ and σ^2 approach either $+\infty$ or $-\infty$ (at the same time) may lead to undefined forms. Sufficiently strong side conditions, such as having the sequences increase in closely linked and particular ways, are needed to create valid distributions. We leave consideration of such exotic cases to others.

Having carefully defined the scalar Gaussian distribution, it would be natural to derive associated properties, such as the characteristic function, MGF, CGF, moments, and cumulants. In particular, $y \sim (\mathcal{S})\mathcal{N}(\mu, \sigma^2)$ has mean $E(y) = \mu$ and variance $\mathcal{V}(y) = \sigma^2$ ($-\infty < \mu < \infty$ and $0 \leq \sigma^2 < \infty$). Given that the reader is probably already familiar with properties of the scalar Gaussian, we present the scalar forms only as special cases of vector forms in the next section.

8.3 THE VECTOR (“MULTIVARIATE”) GAUSSIAN DISTRIBUTION

Our approach to the vector Gaussian proceeds in three stages. First we construct random vectors as increasingly more complex linear transformations of a set of i.i.d. scalar Gaussian variables $\{z_i\} \sim \mathcal{N}(0, 1)$. Second, we describe the associated characteristic and related generating functions. Third, and finally, we demonstrate how any random vector possessing such a characteristic function can be expressed in terms of an underlying set of i.i.d. scalar Gaussian variables. No claim is made that the underlying variables were constructed directly from independent Gaussian variables. However, we are free to operate as though they were. Consequently properties of *any* vector Gaussian can be expressed in terms of properties of a set of fully independent unit Gaussian variables.

Definition 8.2 (a) An $m \times 1$ random vector \mathbf{z} ($1 \leq m < \infty$) with i.i.d. element $z_j \sim \mathcal{N}(0, 1)$ follows a *standard vector (multivariate) Gaussian distribution*, written $\mathbf{z} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$, with $\mathbf{0}$ $m \times 1$ and \mathbf{I} $m \times m$.

(b) An $m \times 1$ random vector \mathbf{z} ($1 \leq m < \infty$) with i.i.d. element $z_j \sim \mathcal{N}(\mu_j, 1)$ follows a *standard noncentral vector (multivariate) Gaussian distribution*, written $\mathbf{z} \sim \mathcal{N}_m(\boldsymbol{\mu}, \mathbf{I})$, with $\boldsymbol{\mu} = \{\mu_j\}$ $m \times 1$ and \mathbf{I} $m \times m$. For clarity, if $\boldsymbol{\mu} = \mathbf{0}$, then \mathbf{z} may be described as standard *central vector (multivariate) Gaussian*.

(c) If $\mathbf{z} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$ and finite constant Φ is $n \times m$ with $1 \leq \text{rank}(\Phi) = m \leq n$, the $n \times 1$ vector $\mathbf{y} = \Phi\mathbf{z} + \boldsymbol{\mu}$ follows a *vector* (“multivariate”) *Gaussian distribution*, written $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$, with $\Sigma = \Phi\Phi'$ ($n \times n$) of rank m .

(d) If $m = n$, then Σ is nonsingular and \mathbf{y} follows a *nonsingular vector* (multivariate) *Gaussian distribution*.

(e) If $m < n$, then Σ is singular and \mathbf{y} follows a *singular vector* (multivariate) *Gaussian distribution*, written $\mathbf{y} \sim S\mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$.

(f) If $m \leq n$, then Σ may or may not be singular, which may be indicated $\mathbf{y} \sim (S)\mathcal{N}_n(\boldsymbol{\mu}_y, \Sigma)$ for clarity, simply as $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}_y, \Sigma)$.

(g) If $\mathbf{z} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$ while finite constant matrix Φ is $n \times m$, $n \geq m \geq 1$, and $\text{rank}(\Phi) = 0$, then $\Phi = \mathbf{0}$ and $n \times 1$ vector $\mathbf{y} = \Phi\mathbf{z} + \boldsymbol{\mu}$, for $\boldsymbol{\mu}$ finite and constant, follows a *degenerate* (or completely singular) *vector* (multivariate) *Gaussian distribution*, written $\mathbf{y} \sim S\mathcal{N}_n(\boldsymbol{\mu}, \mathbf{0})$. Also $\Pr\{\mathbf{y} = \boldsymbol{\mu}\} = 1$.

Example 8.1 A $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors and full rank \mathbf{X} has

$$\begin{aligned}\hat{\boldsymbol{\beta}} &= [(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y} \\ &\sim \mathcal{N}_q[\boldsymbol{\beta}, (\mathbf{X}'\mathbf{X})^{-1}\sigma^2].\end{aligned}\quad (8.12)$$

With or without full-rank \mathbf{X} but with the requirements of full rank of $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ and $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X})$ (which ensures $\boldsymbol{\theta}$ is testable),

$$\begin{aligned}\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 &= \mathbf{C}\tilde{\boldsymbol{\beta}} - \boldsymbol{\theta}_0 \\ &\sim \mathcal{N}_a\{\boldsymbol{\theta} - \boldsymbol{\theta}_0, [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}\sigma^2\}.\end{aligned}\quad (8.13)$$

For $\text{rank}(\mathbf{X}) = r \leq q$, $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, $\text{rank}(\mathbf{H}) = r$, and

$$\begin{aligned}\hat{\mathbf{y}} &= \mathbf{H}\mathbf{y} \\ &\sim S\mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \mathbf{H}\sigma^2).\end{aligned}\quad (8.14)$$

Also

$$\begin{aligned}\hat{\mathbf{e}} &= (\mathbf{I} - \mathbf{H})\mathbf{y} \\ &\sim S\mathcal{N}_N[\mathbf{0}, (\mathbf{I} - \mathbf{H})\sigma^2],\end{aligned}\quad (8.15)$$

with $\text{rank}(\mathbf{I} - \mathbf{H}) = N - r$. In contrast, $\mathbf{e} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}\sigma^2)$.

Theorem 8.2 (a) A standard central vector Gaussian $\mathbf{z} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$ has characteristic function $\phi_{\mathbf{z}}(\mathbf{t}) = \exp(-\mathbf{t}'\mathbf{t}/2)$, moment generating function $m_{\mathbf{z}}(\mathbf{t}) = \exp(\mathbf{t}'\mathbf{t}/2)$, mean $E(\mathbf{z}) = \mathbf{0}$ ($m \times 1$) and covariance $\mathcal{V}(\mathbf{z}) = \mathbf{I}$ ($m \times m$). (b) A standard noncentral vector Gaussian $\mathbf{z} \sim \mathcal{N}_m(\boldsymbol{\mu}, \mathbf{I})$ has characteristic function $\phi_{\mathbf{z}}(\mathbf{t}) = \exp(i\mathbf{t}'\boldsymbol{\mu} - \mathbf{t}'\mathbf{t}/2)$, moment generating function $m_{\mathbf{z}}(\mathbf{t}) = \exp(\mathbf{t}'\boldsymbol{\mu} + \mathbf{t}'\mathbf{t}/2)$, mean $E(\mathbf{z}) = \boldsymbol{\mu}$ ($m \times 1$), and covariance $\mathcal{V}(\mathbf{z}) = \mathbf{I}$ ($m \times m$).

(c) The characteristic function of $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is $\phi_{\mathbf{y}}(\mathbf{t}) = \exp(i\mathbf{t}'\boldsymbol{\mu} - \mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}/2)$, and the moment generating function is $m_{\mathbf{y}}(\mathbf{t}) = \exp(\mathbf{t}'\boldsymbol{\mu} + \mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}/2)$.

(d) A vector Gaussian $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ has a distribution fully characterized by $\{\boldsymbol{\mu}, \boldsymbol{\Sigma}\}$, which are the mean vector and covariance matrix, respectively.

Proof of (a). Due to independence properties,

$$\begin{aligned} \phi_{\mathbf{z}}(\mathbf{t}) &= \prod_{j=1}^m \phi_{z_j}(t_j) = \prod_{j=1}^m \exp(-t_j^2/2) \\ &= \exp\left(-\sum_{j=1}^m t_j^2/2\right) \\ &= \exp(-\mathbf{t}'\mathbf{t}/2). \end{aligned} \tag{8.16}$$

A parallel argument gives $m_{\mathbf{z}}(\mathbf{t})$.

The first moment is $E(\mathbf{z}) = \{E(z_j)\} = \{0\}$. Independence among elements gives $\mathcal{V}(z_j, z_{j'}) = 0$. Combining zero covariance with $E(z_j^2) = 1$ gives $\mathcal{V}(\mathbf{z}) = \mathbf{I}$.

Proof of (b). Left as an exercise.

Proof of (c). Necessarily $\boldsymbol{\Sigma} = \boldsymbol{\Phi}\boldsymbol{\Phi}'$ for $\boldsymbol{\Phi}$ $m \times n$ of rank m , and $\mathbf{y} = \boldsymbol{\Phi}\mathbf{z} + \boldsymbol{\mu}$ is a linear transformation of $\mathbf{z} \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$. Linear transformation properties give

$$\begin{aligned} \phi_{\mathbf{y}}(\mathbf{s}) &= \exp(i\mathbf{s}'\boldsymbol{\mu})\phi_{\mathbf{z}}(\boldsymbol{\Phi}'\mathbf{s}) \\ &= \exp(i\mathbf{s}'\boldsymbol{\mu})\exp[-(\boldsymbol{\Phi}'\mathbf{s})'(\boldsymbol{\Phi}'\mathbf{s})/2] \\ &= \exp(i\mathbf{s}'\boldsymbol{\mu} - \mathbf{s}'\boldsymbol{\Phi}\boldsymbol{\Phi}'\mathbf{s}/2) \\ &= \exp(i\mathbf{s}'\boldsymbol{\mu} - \mathbf{s}'\boldsymbol{\Sigma}\mathbf{s}/2). \end{aligned} \tag{8.17}$$

A parallel approach provides the MGF.

Proof of (d). The characterization by $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is obvious from inspection of the characteristic function, coupled with the unique determination of a distribution by a characteristic function. The moments may be found by differentiating the MGF, once for the mean and twice for the second moment about zero, and setting $\mathbf{t} = \mathbf{0}$. Alternately, $E(\mathbf{y}) = E(\boldsymbol{\Phi}\mathbf{z} + \boldsymbol{\mu}) = \boldsymbol{\Phi}E(\mathbf{z}) + \boldsymbol{\mu} = \boldsymbol{\mu}$ and $\mathcal{V}(\mathbf{y}) = \mathcal{V}(\boldsymbol{\Phi}\mathbf{z} + \boldsymbol{\mu}_y) = \mathcal{V}(\boldsymbol{\Phi}\mathbf{z}) = \boldsymbol{\Phi}\mathcal{V}(\mathbf{z})\boldsymbol{\Phi}' = \boldsymbol{\Phi}\boldsymbol{\Phi}' = \boldsymbol{\Sigma}$. \square

Lemma 8.2 Any full-rank and orthonormal transformation of a standard vector Gaussian is also standard vector Gaussian of possibly smaller dimension. More specifically, if \mathbf{A} is constant, $m \times n$, $1 \leq m \leq n$, $\mathbf{A}\mathbf{A}' = \mathbf{I}_m$ (full row rank and rowwise orthonormal), and $\mathbf{z}_1 \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$, then $\mathbf{z}_2 = \mathbf{A}\mathbf{z}_1 \sim \mathcal{N}_m(\mathbf{0}, \mathbf{I})$.

Proof. As a linear transformation of \mathbf{z}_1 , the characteristic function of $m \times 1$ \mathbf{z}_2 , with $m \times 1$ \mathbf{s} , is

$$\begin{aligned}
 \phi_{\mathbf{z}_2}(\mathbf{s}) &= \phi_{\mathbf{A}\mathbf{z}_1}(\mathbf{s}) = e^{i\mathbf{s}'\mathbf{0}} \cdot \phi_{\mathbf{z}_1}(\mathbf{A}'\mathbf{s}) \\
 &= 1 \cdot \exp[-(\mathbf{A}'\mathbf{s})'\mathbf{A}'\mathbf{s}/2] \\
 &= \exp(-\mathbf{s}'\mathbf{s}/2).
 \end{aligned} \tag{8.18}$$

Recognizing the last expression as the characteristic function of a standard vector Gaussian (of dimension m), coupled with the unique determination of a distribution by a characteristic function, completes the proof. \square

Theorem 8.3 Linear transformations of vector Gaussian variables are Gaussian (the reproductive property). If $\mathbf{y}_2 = \mathbf{C}_1\mathbf{y}_1 + \mathbf{c}_0$ for finite constants \mathbf{C}_1 ($n_2 \times n_1$), \mathbf{c}_0 ($n_2 \times 1$), and $\mathbf{y}_1 \sim \mathcal{N}_{n_1}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_1)$, then $\mathbf{y}_2 \sim \mathcal{N}_{n_2}(\mathbf{C}_1\boldsymbol{\mu}_1 + \mathbf{c}_0, \mathbf{C}_1\boldsymbol{\Sigma}_1\mathbf{C}_1')$. The special case of $n_2 = 1$ implies all linear combinations of vector Gaussian variables are scalar Gaussian. As for all Gaussian variables, the rank of $\boldsymbol{\Sigma}_2 = \mathbf{C}_1\boldsymbol{\Sigma}_1\mathbf{C}_1'$ determines whether \mathbf{y}_2 is degenerate, singular, or nonsingular.

Proof. As a linear transformation, the characteristic function may written

$$\begin{aligned}
 \phi_{\mathbf{y}_2}(\mathbf{s}) &= \exp(i\mathbf{s}'\mathbf{c}_0)\phi_{\mathbf{y}_1}(\mathbf{C}_1'\mathbf{s}) \\
 &= \exp(i\mathbf{s}'\mathbf{c}_0)\exp[i(\mathbf{C}_1'\mathbf{s})'\boldsymbol{\mu}_1 - (\mathbf{C}_1'\mathbf{s})'\boldsymbol{\Sigma}_1(\mathbf{C}_1'\mathbf{s})/2] \\
 &= \exp(i\mathbf{s}'\mathbf{c}_0 + i\mathbf{s}'\mathbf{C}_1\boldsymbol{\mu}_1 - \mathbf{s}'\mathbf{C}_1\boldsymbol{\Sigma}_1\mathbf{C}_1'\mathbf{s}/2) \\
 &= \exp[i\mathbf{s}'(\mathbf{c}_0 + \mathbf{C}_1\boldsymbol{\mu}_1) - \mathbf{s}'\mathbf{C}_1\boldsymbol{\Sigma}_1\mathbf{C}_1'\mathbf{s}/2].
 \end{aligned} \tag{8.19}$$

The function is the characteristic function of $\mathbf{y}_2 \sim \mathcal{N}_{n_2}(\mathbf{C}_1\boldsymbol{\mu}_1 + \mathbf{c}_0, \mathbf{C}_1\boldsymbol{\Sigma}_1\mathbf{C}_1')$. \square

Theorem 8.4 The central standard vector Gaussian $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$ has density $f_{\mathbf{z}}(\mathbf{z}_*) = (2\pi)^{-n/2}\exp(-\mathbf{z}'_*\mathbf{z}_*/2)$.

Proof. Due to independence properties,

$$\begin{aligned}
 f_{\mathbf{z}}(\mathbf{z}_*) &= \prod_{j=1}^n f(z_{*j}) = \prod_{j=1}^n [(2\pi)^{-1/2} \exp(-z_{*j}^2/2)] \\
 &= (2\pi)^{-n/2} \exp\left(-\sum_{j=1}^n z_{*j}^2/2\right) \\
 &= (2\pi)^{-n/2} \exp(-\mathbf{z}'_*\mathbf{z}_*/2).
 \end{aligned} \tag{8.20}$$

\square

Corollary 8.4 If $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$ and $\mathbf{y} = \boldsymbol{\Phi}\mathbf{z} + \boldsymbol{\mu}$, with $\boldsymbol{\Phi}$ $n \times n$ and rank n , then $\boldsymbol{\Sigma} = \boldsymbol{\Phi}\boldsymbol{\Phi}'$ is full rank (of n , nonsingular). The nonsingular vector Gaussian $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is continuous in n dimensions and absolutely continuous, with density

$$f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\mu}, \boldsymbol{\Sigma}) = (2\pi)^{-n/2} |\boldsymbol{\Sigma}|^{-1/2} \exp[-(\mathbf{y}_* - \boldsymbol{\mu})'\boldsymbol{\Sigma}^{-1}(\mathbf{y}_* - \boldsymbol{\mu})/2]. \tag{8.21}$$

Proof. Full rank of the $n \times n$ matrix $\Sigma = \Phi\Phi'$ follows from full rank of $n \times n$ Φ . The spectral decomposition allows writing $\Sigma = \Upsilon \text{Dg}(\lambda) \Upsilon'$, with $\Upsilon' \Upsilon = \Upsilon \Upsilon' = I_n$. Both Υ and $\text{Dg}(\lambda)$ are square ($n \times n$) and full rank, which gives $\Sigma^{-1} = \Upsilon \text{Dg}(\lambda)^{-1} \Upsilon'$ and $\Phi = \Upsilon \text{Dg}(\lambda)^{1/2}$. In turn,

$$\mathbf{y} = \Phi \mathbf{z} + \boldsymbol{\mu}, \tag{8.22}$$

and \mathbf{y} is a full-rank linear transformation of $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$. The theorem ensures it is continuous in n dimensions and absolutely continuous, with density $f_{\mathbf{z}}(\mathbf{z}_*) = (2\pi)^{-n/2} \exp(-\mathbf{z}'_* \mathbf{z}_*/2)$. As a smooth and one-to-one function of an a.c. random vector, \mathbf{y} is also. By Theorem 7.17

$$\begin{aligned} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\mu}, \Sigma) &= \|\Phi^{-1}\| f_{\mathbf{z}}[\Phi^{-1}(\mathbf{y}_* - \boldsymbol{\mu})] \\ &= \|\text{Dg}(\lambda)^{-1/2}\| \|\Upsilon'\| f_{\mathbf{z}_R}[\Phi^{-1}(\mathbf{y}_* - \boldsymbol{\mu})] \\ &= |\Sigma|^{-1/2} (2\pi)^{-n/2} \exp\left\{-[\Phi^{-1}(\mathbf{y}_* - \boldsymbol{\mu})]' [\Phi^{-1}(\mathbf{y}_* - \boldsymbol{\mu})]/2\right\} \\ &= |\Sigma|^{-1/2} (2\pi)^{-n/2} \exp\left[-(\mathbf{y}_* - \boldsymbol{\mu})' [\Phi^{-t} \Phi^{-1}(\mathbf{y}_* - \boldsymbol{\mu})]/2\right] \\ &= |\Sigma|^{-1/2} (2\pi)^{-n/2} \exp\left[-(\mathbf{y}_* - \boldsymbol{\mu})' \Sigma^{-1}(\mathbf{y}_* - \boldsymbol{\mu})/2\right]. \quad \square \end{aligned} \tag{8.23}$$

Theorem 8.5 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$ and $\Sigma = \Phi\Phi'$, with Φ $n \times n_1$ of rank n_1 with $0 < n_1 < n$, then $\mathbf{y} = \Phi \mathbf{z} + \boldsymbol{\mu}$ with $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I})$, $\text{rank}(\Sigma) = n_1$ (nondegenerate but less than full rank, singular). The vector Gaussian \mathbf{y} can also be expressed as the sum of linear transformations of a nonsingular vector Gaussian \mathbf{y}_1 of dimension n_1 (continuous and a.c.) and a discrete and degenerate random vector \mathbf{y}_2 of dimension $n - n_1$. Furthermore, the distribution function of \mathbf{y} can be stated conveniently in terms of the density of the nonsingular vector Gaussian of dimension n_1 . Particular forms for the results may be expressed in terms of the spectral decomposition of the $n \times n$ matrix Σ . With λ_1 the $n_1 \times 1$ vector of positive eigenvalues, $n_2 = n - n_1$, $\mathbf{0}$ of dimension $n_2 \times 1$, columnwise orthonormal Υ_1 of dimension $n \times n_1$, columnwise orthonormal Υ_2 of $n \times n_2$, and $\Upsilon_1' \Upsilon_2 = \mathbf{0}$,

$$\begin{aligned} \Sigma &= \Upsilon \text{Dg}(\lambda_1, \mathbf{0}) \Upsilon' \\ &= [\Upsilon_1 \quad \Upsilon_2] \text{Dg}(\lambda_1, \mathbf{0}) \begin{bmatrix} \Upsilon_1' \\ \Upsilon_2' \end{bmatrix}. \end{aligned} \tag{8.24}$$

If

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \Upsilon_1' \\ \Upsilon_2' \end{bmatrix} \mathbf{y} = \begin{bmatrix} \Upsilon_1' \mathbf{y} \\ \Upsilon_2' \mathbf{y} \end{bmatrix} \tag{8.25}$$

$$\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} = \begin{bmatrix} \Upsilon_1' \\ \Upsilon_2' \end{bmatrix} \boldsymbol{\mu} = \begin{bmatrix} \Upsilon_1' \boldsymbol{\mu} \\ \Upsilon_2' \boldsymbol{\mu} \end{bmatrix}, \tag{8.26}$$

then $\mathbf{y}_1 \sim \mathcal{N}_{n_1}[\boldsymbol{\mu}_1, \text{Dg}(\lambda_1)]$ is a nonsingular vector Gaussian with density $f_{\mathbf{y}_1}[\mathbf{y}_{*1}; \boldsymbol{\mu}_1, \text{Dg}(\lambda_1)]$ and $\mathbf{y}_2 \sim \mathcal{SN}_{n_2}(\boldsymbol{\mu}_2, \mathbf{0})$ is discrete and degenerate. Finally,

the distribution function of \mathbf{y} can be conveniently expressed in terms of $K_g = \prod_{j=1}^{n_1} (2\pi\lambda_j)^{-1/2}$ and the function

$$\begin{aligned} g(\mathbf{y}_*; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= K_g \exp[-(\mathbf{y}_* - \boldsymbol{\mu})' \boldsymbol{\Sigma}^+ (\mathbf{y}_* - \boldsymbol{\mu})/2] \\ &= K_g \exp\left[-(\mathbf{y}_* - \boldsymbol{\mu})' \begin{bmatrix} \boldsymbol{\Upsilon}_1 & \boldsymbol{\Upsilon}_2 \end{bmatrix} \text{Dg}(\boldsymbol{\lambda}_1, \mathbf{0})^+ \begin{bmatrix} \boldsymbol{\Upsilon}'_1 \\ \boldsymbol{\Upsilon}'_2 \end{bmatrix} (\mathbf{y}_* - \boldsymbol{\mu})/2\right] \\ &= K_g \exp\left[-\begin{bmatrix} (\mathbf{y}_{*1} - \boldsymbol{\mu}_1) \\ (\mathbf{y}_{*2} - \boldsymbol{\mu}_2) \end{bmatrix}' \begin{bmatrix} \text{Dg}(\boldsymbol{\lambda}_1)^{-1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} (\mathbf{y}_{*1} - \boldsymbol{\mu}_1) \\ (\mathbf{y}_{*2} - \boldsymbol{\mu}_2) \end{bmatrix} /2\right] \\ &= f_{\mathbf{y}_1}[\mathbf{y}_{*1}; \boldsymbol{\mu}_1, \text{Dg}(\boldsymbol{\lambda}_1)]. \end{aligned} \quad (8.27)$$

Function $g(\mathbf{y}_*; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ serves as a pseudodensity for \mathbf{y} in the following sense. With $n_1 \times 1$ $\mathbf{y}_{*1} = \boldsymbol{\Upsilon}'_1 \mathbf{y}_*$,

$$\begin{aligned} \Pr\{\mathbf{y} \leq \mathbf{y}_0\} &= \Pr\{\mathbf{y}_1 \leq \mathbf{y}_{0,1}\} \\ &= \int_{-\infty}^{y_{0,1,1}} \cdots \int_{-\infty}^{y_{0,1,n_1}} f_{\mathbf{y}_1}[\mathbf{y}_{*1}; \boldsymbol{\mu}_1, \text{Dg}(\boldsymbol{\lambda}_1)] dy_{*1,1} \cdots dy_{*1,n_1}. \end{aligned} \quad (8.28)$$

However, the integration is over $\mathfrak{R}^{n_1} \subset \mathfrak{R}^n$ (and therefore does not meet all requirements for the definition of a density). Explicit and valid integration over \mathfrak{R}^n requires a more general type of integration (Riemann-Stieltjes).

Proof. The outer product $\boldsymbol{\Sigma}$ is symmetric and either positive definite or positive semidefinite. Having $\text{rank}(\boldsymbol{\Sigma}) = n_1$ and $0 < n_1 < n$ ensures $n_0 = n - n_1 > 0$. Therefore the dimensions and nature of the spectral decomposition are as claimed.

Expressing \mathbf{y} in terms of a sum of nonsingular and degenerate Gaussian variables is achieved simply by considering three particular characteristic functions. Applying the theorem for the characteristic function of a linear transformation to $\mathbf{y}_1 = \boldsymbol{\Upsilon}'_1 \mathbf{y}$ gives $\phi_{\mathbf{y}_1}(\mathbf{t}_1) = \exp[i\mathbf{t}'_1 \boldsymbol{\mu}_1 - \mathbf{t}'_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{t}_1/2]$ as the characteristic function of the $n_1 \times 1$ random vector $\mathbf{y}_1 \sim \mathcal{N}_{n_1}[\boldsymbol{\mu}_1, \text{Dg}(\boldsymbol{\lambda}_1)]$. The full rank of n_1 for $\text{Dg}(\boldsymbol{\lambda}_1)$ ensures \mathbf{y}_1 is continuous and a.c. (by the previous theorem). Similar reasoning with $\mathbf{y}_2 = \boldsymbol{\Upsilon}'_2 \mathbf{y}$ gives $\phi_{\mathbf{y}_2}(\mathbf{t}_2) = \exp(i\mathbf{t}'_2 \boldsymbol{\mu}_2)$ as the characteristic function of the $n_2 \times 1$ degenerate random vector $\mathbf{y}_2 \sim \mathcal{SN}_{n_2}(\boldsymbol{\mu}_2, \mathbf{0})$, with $\Pr\{\mathbf{y}_2 = \boldsymbol{\mu}_2\} = 1$ (and therefore discrete). Using the spectral decomposition of $\boldsymbol{\Sigma}$ to expand and then simplify the characteristic function of \mathbf{y} gives

$$\begin{aligned} \phi_{\mathbf{y}}(\mathbf{t}) &= \exp(i\mathbf{t}' \boldsymbol{\mu} - \mathbf{t}' \boldsymbol{\Sigma} \mathbf{t}/2) \\ &= \exp\left(i\mathbf{t}' \begin{bmatrix} \boldsymbol{\Upsilon}_1 & \boldsymbol{\Upsilon}_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\Upsilon}'_1 \\ \boldsymbol{\Upsilon}'_2 \end{bmatrix} \boldsymbol{\mu} - \mathbf{t}' \begin{bmatrix} \boldsymbol{\Upsilon}_1 & \boldsymbol{\Upsilon}_2 \end{bmatrix} \text{Dg}(\boldsymbol{\lambda}_1, \mathbf{0}) \begin{bmatrix} \boldsymbol{\Upsilon}'_1 \\ \boldsymbol{\Upsilon}'_2 \end{bmatrix} \mathbf{t}/2\right) \\ &= \exp\left(i \begin{bmatrix} \mathbf{t}'_1 & \mathbf{t}'_2 \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} - \begin{bmatrix} \mathbf{t}'_1 & \mathbf{t}'_2 \end{bmatrix} \text{Dg}(\boldsymbol{\lambda}_1, \mathbf{0}) \begin{bmatrix} \mathbf{t}_1 \\ \mathbf{t}_2 \end{bmatrix} /2\right) \\ &= \exp[i(\mathbf{t}'_2 \boldsymbol{\mu}_2 + \mathbf{t}'_1 \boldsymbol{\mu}_1) - \mathbf{t}'_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{t}_1/2] \\ &= \exp(i\mathbf{t}'_2 \boldsymbol{\mu}_2) \exp[i\mathbf{t}'_1 \boldsymbol{\mu}_1 - \mathbf{t}'_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{t}_1/2] \\ &= \phi_{\mathbf{y}_2}(\mathbf{t}_2) \phi_{\mathbf{y}_1}(\mathbf{t}_1). \end{aligned} \quad (8.29)$$

The last form implies \mathbf{y}_1 and \mathbf{y}_2 are statistically independent. As eigenvectors of a

symmetric matrix (chosen to be orthonormal without loss of generality), $\Upsilon = [\Upsilon_1 \ \Upsilon_2]$ is square, full rank, $\Upsilon'\Upsilon = \Upsilon\Upsilon'$, and $(\Upsilon)^{-1} = \Upsilon'$. Hence

$$\begin{aligned} \mathbf{y} &= [\Upsilon_1 \ \Upsilon_2] \begin{bmatrix} \Upsilon'_1 \\ \Upsilon'_2 \end{bmatrix} \mathbf{y} \\ &= [\Upsilon_1 \ \Upsilon_2] \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \Upsilon_1 \mathbf{y}_1 + \Upsilon_2 \mathbf{y}_2, \end{aligned} \tag{8.30}$$

which demonstrates \mathbf{y} is a full-rank (and invertible) linear transformation of $[\mathbf{y}'_1 \ \mathbf{y}'_2]$ and also a sum of linear transformations of \mathbf{y}_1 and \mathbf{y}_2 .

If $g(\mathbf{y}_*; \boldsymbol{\mu}, \Sigma)$ allows computing the CDF of \mathbf{y} as claimed, then it must be the density of \mathbf{y}_1 . If \mathbf{A} is $m \times n$, \mathbf{P} is $h \times m$ with $\mathbf{P}\mathbf{P} = \mathbf{I}_m$, and \mathbf{Q} is $n \times p$ with $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_p$, then $(\mathbf{P}\mathbf{A}\mathbf{Q})^+ = \mathbf{Q}^+ \mathbf{A}^+ \mathbf{P}^+ = \mathbf{Q}' \mathbf{A}^+ \mathbf{P}'$ (Lemma 1.15). Therefore $\Sigma^+ = \Upsilon \text{Dg}(\lambda_1, \mathbf{0})^+ \Upsilon'$. With $K = (2\pi)^{-n_1/2} \prod_{j=1}^{n_1} \lambda_j^{-1/2}$, the function is

$$\begin{aligned} f[\mathbf{y}_{*1}; \boldsymbol{\mu}_1, \text{Dg}(\lambda_1)] &= K \exp[-(\mathbf{y}_{*1} - \boldsymbol{\mu}_1)' \text{Dg}(\lambda_1)^{-1} (\mathbf{y}_{*1} - \boldsymbol{\mu}_1)/2] \\ &= K \exp \left[- \begin{bmatrix} (\mathbf{y}_{*1} - \boldsymbol{\mu}_1) \\ (\mathbf{y}_{*2} - \boldsymbol{\mu}_2) \end{bmatrix}' \begin{bmatrix} \text{Dg}(\lambda_1)^{-1} \mathbf{0} \\ \mathbf{0} \ \mathbf{0} \end{bmatrix} \begin{bmatrix} (\mathbf{y}_{*1} - \boldsymbol{\mu}_1) \\ (\mathbf{y}_{*2} - \boldsymbol{\mu}_2) \end{bmatrix} / 2 \right] \\ &= K \exp \left[- \begin{bmatrix} (\Upsilon'_1 \mathbf{y}_{*1} - \Upsilon'_1 \boldsymbol{\mu}) \\ (\Upsilon'_2 \mathbf{y}_{*2} - \Upsilon'_2 \boldsymbol{\mu}) \end{bmatrix}' \begin{bmatrix} \text{Dg}(\lambda_1)^{-1} \mathbf{0} \\ \mathbf{0} \ \mathbf{0} \end{bmatrix} \begin{bmatrix} (\Upsilon'_1 \mathbf{y}_{*1} - \Upsilon'_1 \boldsymbol{\mu}) \\ (\Upsilon'_2 \mathbf{y}_{*2} - \Upsilon'_2 \boldsymbol{\mu}) \end{bmatrix} / 2 \right] \\ &= K \exp \left[-(\mathbf{y}_* - \boldsymbol{\mu})' [\Upsilon_1 \ \Upsilon_2] \text{Dg}(\lambda_1, \mathbf{0})^+ \begin{bmatrix} \Upsilon'_1 \\ \Upsilon'_2 \end{bmatrix} (\mathbf{y}_* - \boldsymbol{\mu}) / 2 \right] \\ &= K \exp[-(\mathbf{y}_* - \boldsymbol{\mu})' \Sigma^+ (\mathbf{y}_* - \boldsymbol{\mu}) / 2] \\ &= g(\mathbf{y}_*; \boldsymbol{\mu}, \Sigma). \end{aligned} \tag{8.31}$$

If $g(\mathbf{y}_*; \boldsymbol{\mu}, \Sigma)$ allows computing the CDF of \mathbf{y} , it must suffice to consider only \mathbf{y}_1 in computing the CDF of \mathbf{y} . Here \mathbf{y} and $[\mathbf{y}'_1 \ \mathbf{y}'_2]$ are one-to-one functions of each other because $\mathbf{y} = \Upsilon [\mathbf{y}'_1 \ \mathbf{y}'_2]'$, although the transformation is not smooth. Both are of dimension $n = n_1 + n_2 > n_1$, while \mathbf{y}_1 has dimension n_1 and \mathbf{y}_2 has dimension n_2 . The degenerate nature of \mathbf{y}_2 gives $\Pr\{\mathbf{y}_2 = \boldsymbol{\mu}_0\} = 1$, which ensures $[\mathbf{y}'_1 \ \mathbf{y}'_2] = [\mathbf{y}'_1 \ \boldsymbol{\mu}'_2]$, with probability 1. Therefore the n_1 dimensions of \mathbf{y}_1 capture all of the randomness in \mathbf{y} , while the n_2 dimensions of \mathbf{y}_2 affect only the location (the mean of \mathbf{y}). Intuitively, a constant carries no information about a (truly) random vector.

The argument may be formalized precisely by generalizing the concept of integration to cover computing probabilities for any combination of continuous and discrete random vectors. The interested reader may consult Lindgren (1976) for a more detailed treatment of the basic ideas behind Riemann-Stieltjes integration. \square

The dimensions of the matrices change between the first two equations in the sequence, corresponding to a change from \mathfrak{R}^{n_1} to $\mathfrak{R}^n = \mathfrak{R}^{n_1} + \mathfrak{R}^{n_0}$. Here the plus sign indicates no overlap between the two subspaces. The form $\mathbf{y}_2 - \boldsymbol{\mu}_2$ occurs only with the additional n_2 -dimensional subspace and is guaranteed to be

$\boldsymbol{\mu}_2 - \boldsymbol{\mu}_2 = \mathbf{0}$. The expression for $g()$ is not unique because the same result may be computed by replacing Σ^+ with any choice of Σ^- (of which infinitely many exist, except when Σ is full rank).

It is often convenient to express $\mathbf{y} \sim (\mathcal{S})\mathcal{N}_n(\boldsymbol{\mu}_y, \Sigma)$, with $\Sigma = \Phi\Phi'$, explicitly in terms of i.i.d. standard Gaussian variables. If $\text{rank}(\Sigma) = n$, then Φ is $n \times n$ with full rank and $\mathbf{z} = \Phi^{-1}(\mathbf{y} - \boldsymbol{\mu}_y) \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I})$. In turn, $\mathbf{y} = \Phi(\mathbf{z} + \boldsymbol{\mu}_z)$ with $\boldsymbol{\mu}_z = \Phi^{-1}\boldsymbol{\mu}_y$ and $(\mathbf{z} + \boldsymbol{\mu}_z) \sim \mathcal{N}_n(\boldsymbol{\mu}_z, \mathbf{I})$. The following lemma generalizes the decomposition to the singular case. The result is very helpful in deriving the distribution of a general quadratic form, as will be seen in the next chapter.

Lemma 8.3 If $\mathbf{y} \sim (\mathcal{S})\mathcal{N}_n(\boldsymbol{\mu}_y, \Sigma)$ with finite $\boldsymbol{\mu}_y$ and finite Σ of rank $0 < n_1 \leq n$, then $n \times n_1$ finite constant Φ_1 exists such that $\Sigma = \Phi_1\Phi_1'$. If $\boldsymbol{\mu}_z = (\Phi_1'\Phi_1)^{-1}\Phi_1'\boldsymbol{\mu}_y = (\Phi_1^+)' \boldsymbol{\mu}_y$ and $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I}_{n_1})$, then $\mathbf{y} = \Phi_1(\mathbf{z} + \boldsymbol{\mu}_z)$.

Proof. Spectral decomposition gives $\Sigma = \Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1'$ with $\Upsilon_1' \Upsilon_1 = \mathbf{I}_{n_1}$ and $\Phi_1 = \Upsilon_1 \text{Dg}(\lambda_1)^{1/2}$. In turn $\Phi_1^+ = \Phi_1'(\Phi_1'\Phi_1)^{-1} = \text{Dg}(\lambda_1)^{-1/2} \Upsilon_1'$. Furthermore $\mathcal{V}(\mathbf{y}) = \Phi_1 \mathbf{I}_{n_1} \Phi_1' = \Sigma$. Requiring $\boldsymbol{\mu}_y = \Phi_1 \boldsymbol{\mu}_z$ implies $(\Phi_1'\Phi_1)^{-1}\Phi_1'\boldsymbol{\mu}_y = \boldsymbol{\mu}_z$. As a linear transformation of \mathbf{z} , the vector \mathbf{y} is Gaussian. \square

8.4 MARGINAL DISTRIBUTIONS

Theorem 8.6 If $\mathbf{y} \sim \mathcal{N}_{n_1+n_2}(\boldsymbol{\mu}, \Sigma)$ with

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \quad (8.32)$$

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} \quad (8.33)$$

$$\Sigma = \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix}, \quad (8.34)$$

\mathbf{y}_j and $\boldsymbol{\mu}_j$ are $n_j \times 1$ for $j \in \{1, 2\}$, while Σ_{jj} is $n_j \times n_j$, then the marginal distribution of \mathbf{y}_j is $\mathcal{N}_{n_j}(\boldsymbol{\mu}_j, \Sigma_{jj})$. More generally, all marginal distributions of a vector Gaussian are vector Gaussian.

Proof. The characteristic function of the marginal of \mathbf{y}_1 is

$$\begin{aligned} \phi_{\mathbf{y}_1}(\mathbf{t}_1) &= \phi_{\mathbf{y}}\left(\begin{bmatrix} \mathbf{t}_1 \\ \mathbf{0} \end{bmatrix}\right) \\ &= \exp\left\{i[\mathbf{t}_1' \ \mathbf{0}'] \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix} - [\mathbf{t}_1' \ \mathbf{0}'] \begin{bmatrix} \Sigma_{11} & \Sigma_{21} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \begin{bmatrix} \mathbf{t}_1 \\ \mathbf{0} \end{bmatrix} / 2\right\} \\ &= \exp(i\mathbf{t}_1' \boldsymbol{\mu}_1 - \mathbf{t}_1' \Sigma_{11} \mathbf{t}_1 / 2), \end{aligned} \quad (8.35)$$

which is the characteristic function of the $\mathcal{N}_{n_1}(\boldsymbol{\mu}_1, \Sigma_{11})$ distribution. More generally, the result does not depend on the order of the variables. It applies to any rearrangement and therefore to all subsets of $\{y_i\}$. \square

Finding a marginal distribution of a vector Gaussian is very easy. One simply selects the corresponding elements of $\boldsymbol{\mu}$ and corresponding rows and columns of $\boldsymbol{\Sigma}$ to get the parameters of the marginal Gaussian. The property does not necessarily hold for other distributions.

The converse of the theorem is not true, because joint distributions do exist which have Gaussian marginal distributions but do *not* have a *joint* Gaussian distribution. Having Gaussian marginal distributions does *not* ensure a Gaussian joint distribution. More specifically, having $\mathbf{y}_j \sim \mathcal{N}_{n_j}(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_{jj})$ for $j \in \{1, 2\}$ and $\boldsymbol{\Sigma} = \mathcal{V}(\mathbf{y}) = \mathcal{V}\left(\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix}\right) = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{21} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}$ does *not* guarantee $\mathbf{y} \sim \mathcal{N}_{n_1+n_2}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$. Proving the result may be a useful exercise.

8.5 INDEPENDENCE

Theorem 8.7 Assuming $\mathbf{y} \sim \mathcal{N}_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with \mathbf{y} , $\boldsymbol{\mu}$, and $\boldsymbol{\Sigma}$ partitioned correspondingly allows writing

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_J \end{bmatrix} \sim \mathcal{N}_N\left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \\ \vdots \\ \boldsymbol{\mu}_J \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} & \cdots & \boldsymbol{\Sigma}_{1J} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} & \cdots & \boldsymbol{\Sigma}_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ \boldsymbol{\Sigma}_{J1} & \boldsymbol{\Sigma}_{J2} & \cdots & \boldsymbol{\Sigma}_{JJ} \end{bmatrix}\right), \quad (8.36)$$

with \mathbf{y}_j of dimension $n_j \times 1$ and $\sum_{j=1}^J n_j = N$.

The J random vectors $\{\mathbf{y}_j, j \in \{1, 2, \dots, k\}\}$ are mutually (totally) independent

(a) if and only if \mathbf{y}_j and $\mathbf{y}_{j'}$ are (pairwise) independent for all $j \neq j'$ and

(b) if and only if $\boldsymbol{\Sigma}_{jj'} = \mathbf{0}$ for all $j \neq j'$.

(c) The result does *not* hold for all distributions that are not Gaussian.

Proof. The characteristic function for $\mathbf{y} \sim \mathcal{N}_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is defined $\forall \mathbf{t} \in \mathfrak{R}^N$ as $\phi_{\mathbf{y}}(\mathbf{t}) = \exp(i\mathbf{t}'\boldsymbol{\mu} - \mathbf{t}'\boldsymbol{\Sigma}\mathbf{t}/2)$. Partitioning \mathbf{t} to match \mathbf{y} gives $\mathbf{t}' = [\mathbf{t}'_1 \ \mathbf{t}'_2 \ \cdots \ \mathbf{t}'_J]$. Therefore

$$\begin{aligned} \phi_{\mathbf{y}}(\mathbf{t}) &= \exp\left(i\sum_{j=1}^J \mathbf{t}'_j \boldsymbol{\mu}_j - \sum_{j=1}^J \sum_{j'=1}^J \mathbf{t}'_j \boldsymbol{\Sigma}_{jj'} \mathbf{t}_{j'}/2\right) \\ &= \exp\left(i\sum_{j=1}^J \mathbf{t}'_j \boldsymbol{\mu}_j - \sum_{j=1}^J \mathbf{t}'_j \boldsymbol{\Sigma}_{jj} \mathbf{t}_j/2 - \sum_{j=1}^{J-1} \sum_{j'=j+1}^J \mathbf{t}'_j \boldsymbol{\Sigma}_{jj'} \mathbf{t}_{j'}\right). \end{aligned} \quad (8.37)$$

The third summation is zero $\forall \mathbf{t} \in \mathfrak{R}^N \Leftrightarrow \boldsymbol{\Sigma}_{jj'} = \mathbf{0} \ \forall j \neq j'$. Thus $\boldsymbol{\Sigma}_{jj'} = \mathbf{0} \ \forall j \neq j'$ if and only if

$$\begin{aligned}
\phi_{\mathbf{y}}(\mathbf{t}) &= \exp\left(i\sum_{j=1}^J \mathbf{t}'_j \boldsymbol{\mu}_j - \sum_{j=1}^J \mathbf{t}'_j \boldsymbol{\Sigma}_{jj} \mathbf{t}_j / 2\right) \\
&= \prod_{j=1}^J \exp(i\mathbf{t}'_j \boldsymbol{\mu}_j - \mathbf{t}'_j \boldsymbol{\Sigma}_{jj} \mathbf{t}_j / 2) \\
&= \prod_{j=1}^J \phi_j(\mathbf{t}_j).
\end{aligned} \tag{8.38}$$

Mutual independence follows from the factorization. Gaussian distribution of each marginal follows from the form of the individual CFs. \square

Corollary 8.7.1 If $n_j \equiv 1$ and $\mathbf{y} \sim \mathcal{N}_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\boldsymbol{\Sigma} = \text{Dg}(\{\sigma_{11}, \dots, \sigma_{NN}\})$, then the N elements of \mathbf{y} are all mutually independent random variables, and $y_j \sim \mathcal{N}_1(\mu_j, \sigma_{jj}) \forall j$.

Corollary 8.7.2 If $N = 2$, $n_j \geq 1$, and

$$\mathbf{y} = \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} \sim \mathcal{N}_N \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{21} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix} \right), \tag{8.39}$$

then \mathbf{y}_1 and \mathbf{y}_2 are independent if and only if $\boldsymbol{\Sigma}_{12} \equiv \boldsymbol{\Sigma}'_{21} = \mathbf{0}$.

8.6 CONDITIONAL DISTRIBUTIONS

Theorem 8.8 For \mathbf{y}_j and $\boldsymbol{\mu}_j$ of dimension $n_j \times 1$, $\mathbf{y} = [\mathbf{y}'_1 \ \mathbf{y}'_2]'$, $\boldsymbol{\mu} = [\boldsymbol{\mu}'_1 \ \boldsymbol{\mu}'_2]'$, while $\boldsymbol{\Sigma}_{jj'}$ is of dimension $n_j \times n_{j'}$ with

$$\boldsymbol{\Sigma} = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \boldsymbol{\Sigma}_{22} \end{bmatrix}, \tag{8.40}$$

and $N = n_1 + n_2$. Here $\mathbf{y} \sim (\mathcal{S})\mathcal{N}_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\text{rank}(\boldsymbol{\Sigma}) = n_+ \leq N$.

(a) The conditional distribution of $(\mathbf{y}_1 | \mathbf{y}_2 = \mathbf{y}_{02})$ is $\mathcal{N}_{n_1}(\boldsymbol{\mu}_{1.2}, \boldsymbol{\Sigma}_{1.2})$ with $\boldsymbol{\mu}_{1.2} = \boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- (\mathbf{y}_{02} - \boldsymbol{\mu}_2)$ and $\boldsymbol{\Sigma}_{1.2} = \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- \boldsymbol{\Sigma}_{21}$.

(b) Although $\boldsymbol{\Sigma}_{22}^-$ is any generalized inverse of $\boldsymbol{\Sigma}_{22}$, both $\boldsymbol{\mu}_{1.2}$ and $\boldsymbol{\Sigma}_{1.2}$ are invariant to the choice of $\boldsymbol{\Sigma}_{22}^-$ (and therefore can be taken to be $\boldsymbol{\Sigma}_{22}^+$).

(c) $(\mathbf{y}_1 | \mathbf{y}_2 = \mathbf{y}_{02}) \sim (\mathcal{S})\mathcal{N}_{n_1}[\boldsymbol{\mu}_1 + \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^+ (\mathbf{y}_{02} - \boldsymbol{\mu}_2), \boldsymbol{\Sigma}_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^+ \boldsymbol{\Sigma}_{21}]$. (8.41)

Proof. The following are true. (1) A matrix \mathbf{B} exists satisfying $\boldsymbol{\Sigma}_{12} = \mathbf{B}\boldsymbol{\Sigma}_{22}$.

(2) $\mathbf{B} = \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^-$ is a solution, and is unique if and only if $\boldsymbol{\Sigma}_{22}$ is nonsingular.

(3) $\boldsymbol{\Sigma}_{12} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^- \boldsymbol{\Sigma}_{22} = \mathbf{0}$.

Proofs of 1 and 2 are left as an exercise. They are consequences of a slight generalization of the Cauchy-Schwartz inequality for inner products based on positive semidefinite matrices.

The proof of 3 starts with the definition of the generalized inverse, $\Sigma_{22} = \Sigma_{22}\Sigma_{22}^-\Sigma_{22}$. Substituting the expression into the equation $\Sigma_{12} = B\Sigma_{22}$ implies $\Sigma_{12} = B(\Sigma_{22}\Sigma_{22}^-\Sigma_{22}) = (B\Sigma_{22})\Sigma_{22}^-\Sigma_{22} = \Sigma_{12}\Sigma_{22}^-\Sigma_{22}$.

An important result relates stochastic independence and conditional distribution. If \mathbf{y}_1 and \mathbf{y}_2 are random vectors with a joint distribution, then the conditional distribution of $\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02}$, if it exists, is identical to the marginal distribution of \mathbf{y}_1 if and only if \mathbf{y}_1 and \mathbf{y}_2 are independent. The basic approach is to construct a random variable $\mathbf{x} = (\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02})$ using the three results just stated. If $\mathbf{x} = \mathbf{A}(\mathbf{y} - \boldsymbol{\mu}) + \mathbf{c}$ in which

$$\mathbf{A} = \begin{bmatrix} \mathbf{I}_{n_1} & -\Sigma_{12}\Sigma_{22}^- \\ \mathbf{0} & \mathbf{I}_{n_2} \end{bmatrix} \tag{8.42}$$

$$\mathbf{c} = \begin{bmatrix} \boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^-(\mathbf{y}_{02} - \boldsymbol{\mu}_2) \\ \boldsymbol{\mu}_2 \end{bmatrix}, \tag{8.43}$$

then $\mathbf{x} = [\mathbf{x}'_1 \ \mathbf{x}'_2]'$ = \mathbf{y} - adjustment and

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{y}_1 + \Sigma_{12}\Sigma_{22}^-(\mathbf{y}_{02} - \mathbf{y}_2) \\ \mathbf{y}_2 \end{bmatrix}. \tag{8.44}$$

Given $\mathbf{y}_2 = \mathbf{y}_{02}$, the adjustment is zero and $\mathbf{x} = \mathbf{y}$. Furthermore, if $\mathbf{y}_2 = \mathbf{y}_{02}$ then $\mathbf{x}_1 = \mathbf{y}_1$ and $(\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02}) = (\mathbf{x}_1|\mathbf{x}_2 = \mathbf{y}_{02})$. Here $N \times N$ \mathbf{A} is full rank of N and $\text{rank}(\mathbf{A}\Sigma\mathbf{A}') = \text{rank}(\Sigma) = n_+$. As a linear transformation of a vector Gaussian, \mathbf{x} is also vector Gaussian with

$$\begin{aligned} \mathbf{x} &= \mathbf{A}(\mathbf{y} - \boldsymbol{\mu}) + \mathbf{c} \\ &\sim (\mathcal{S})\mathcal{N}_N(\mathbf{c}, \mathbf{A}\Sigma\mathbf{A}'). \end{aligned} \tag{8.45}$$

Using 3 we can prove

$$\mathbf{A}\Sigma\mathbf{A}' = \begin{bmatrix} \Sigma_{11} - \Sigma_{12}\Sigma_{22}^-\Sigma_{21} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{bmatrix}. \tag{8.46}$$

Therefore $\mathbf{x}_1 \perp\!\!\!\perp \mathbf{x}_2$, which ensures the distribution of $(\mathbf{x}_1|\mathbf{x}_2 = \mathbf{y}_{02})$ is identical to the distribution of \mathbf{x}_1 . In turn, the distribution of $(\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02})$ coincides with the distributions of $(\mathbf{x}_1|\mathbf{x}_2 = \mathbf{y}_{02})$ and therefore \mathbf{x}_1 . Finally, $\mathbf{c} = [\boldsymbol{\mu}'_{1.2} \ \boldsymbol{\mu}'_2]'$ and

$$\mathbf{A}\Sigma\mathbf{A}' = \begin{bmatrix} \Sigma_{1.2} & \mathbf{0} \\ \mathbf{0} & \Sigma_{22} \end{bmatrix} \tag{8.47}$$

means $\mathbf{x}_1 \sim (\mathcal{S})\mathcal{N}_{n_1}(\boldsymbol{\mu}_{1.2}, \Sigma_{1.2})$, which is also the distribution of $(\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02})$. \square

Corollary 8.8.1 If Σ is full rank ($n_+ = n_1 + n_2 = N$), then

- (a) Σ_{22} has full rank of n_2 ,
- (b) $\Sigma_{1.2}$ has full rank of n_1 , and
- (c) $(\mathbf{y}_1|\mathbf{y}_2 = \mathbf{y}_{02}) \sim \mathcal{N}_{n_1}(\boldsymbol{\mu}_{1.2}, \Sigma_{1.2})$ has a density with $\boldsymbol{\mu}_{1.2} = \boldsymbol{\mu}_1 + \Sigma_{12}\Sigma_{22}^-(\mathbf{y}_{02} - \boldsymbol{\mu}_2)$ and $\Sigma_{1.2} = \Sigma_{11} - \Sigma_{12}\Sigma_{22}^-\Sigma_{21}$.
- (d) The dispersion matrix $\Sigma_{1.2}$ is not a function of \mathbf{y}_2 , and equals Σ_{11} less an

adjustment which vanishes when $\Sigma_{12} = \mathbf{0}$.

(e) In contrast, mean $\mu_{1.2}$ is a linear function of \mathbf{y}_{02} and equals μ_1 plus an adjustment which vanishes if $\Sigma_{12} = \mathbf{0}$.

Proof. Both \mathbf{y}_1 and \mathbf{y}_2 have marginal densities as well as a joint density. The general form for the PDF of a conditional vector gives, with $\mathbf{y}_{*,0} = [\mathbf{y}'_{*1} \ \mathbf{y}'_{02}]'$,

$$\begin{aligned}
 f_{1|2}(\mathbf{y}_{*1}|\mathbf{y}_2 = \mathbf{y}_{02}) &= \frac{f_{\mathbf{y}}(\mathbf{y}_{*,0}; \boldsymbol{\mu}, \boldsymbol{\Sigma})}{f_{\mathbf{y}_2}(\mathbf{y}_{02}; \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_{22})} \\
 &= \frac{(2\pi)^{-N/2} |\boldsymbol{\Sigma}|^{-1/2} \exp[-(\mathbf{y}_{*,0} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{y}_{*,0} - \boldsymbol{\mu})/2]}{(2\pi)^{-n_2/2} |\boldsymbol{\Sigma}_{22}|^{-1/2} \exp[-(\mathbf{y}_{02} - \boldsymbol{\mu}_2)' \boldsymbol{\Sigma}_{22}^{-1} (\mathbf{y}_{02} - \boldsymbol{\mu}_2)/2]} \\
 &\quad \vdots \quad (\text{details omitted}) \\
 &= (2\pi)^{-n_1/2} |\boldsymbol{\Sigma}_{1.2}|^{-1/2} \exp[-(\mathbf{y}_{*1} - \boldsymbol{\mu}_{1.2})' \boldsymbol{\Sigma}_{1.2}^{-1} (\mathbf{y}_{*1} - \boldsymbol{\mu}_{1.2})/2]. \tag{8.48}
 \end{aligned}$$

Details are left as an exercise. □

Corollary 8.8.2 The validity of a multivariate general linear model data analysis with (random) Gaussian predictors depends on conditional distribution results. For independent sampling unit $i \in \{1, \dots, N\}$, the results can be formally stated in terms of $\mathbf{u}_i = [\mathbf{u}'_{i1} \ \mathbf{u}'_{i2}]' \sim \mathcal{N}_{p+q}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with full rank $\boldsymbol{\Sigma}$ and \mathbf{u}_i independent of $\mathbf{u}_{i'}$ for $i \neq i'$. In turn, $\text{row}_i(\mathbf{Y}) = (\mathbf{u}_{i1} \ \mathbf{u}_{i2} - \boldsymbol{\mu}_{i02})'$ implies \mathbf{Y} is $N \times p$, $\text{row}_i(\mathbf{X}) = [1 \ (\mathbf{u}_{i2} - \boldsymbol{\mu}_{i2})']$ implies \mathbf{X} is $N \times q$, and \mathbf{B} ($q \times p$) is

$$\mathbf{B} = \begin{bmatrix} \boldsymbol{\mu}'_1 \\ (\boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1})' \end{bmatrix}. \tag{8.49}$$

Under the conditions just stated, \mathbf{Y} , \mathbf{X} , and \mathbf{B} satisfy the assumptions of the GLM $_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma}_{1.2})$ with Gaussian errors and

$$\begin{aligned}
 \mathbf{Y} &= \mathbf{X} \mathbf{B} && + \mathbf{E} \\
 \begin{bmatrix} (\mathbf{u}_{11} | \mathbf{u}_{12} = \mathbf{u}_{102})' \\ (\mathbf{u}_{21} | \mathbf{u}_{22} = \mathbf{u}_{202})' \\ \vdots \\ (\mathbf{u}_{N1} | \mathbf{u}_{N2} = \mathbf{u}_{N02})' \end{bmatrix} &= \begin{bmatrix} 1 & (\mathbf{u}_{12} - \boldsymbol{\mu}_{12})' \\ 1 & (\mathbf{u}_{22} - \boldsymbol{\mu}_{22})' \\ \vdots & \vdots \\ 1 & (\mathbf{u}_{N2} - \boldsymbol{\mu}_{N2})' \end{bmatrix} \mathbf{B} + \mathbf{E}. \tag{8.50}
 \end{aligned}$$

Also $E[\text{row}_i(\mathbf{Y})|\mathbf{X}] = \text{row}_i(\mathbf{X})\mathbf{B}$. If $p = 1$, we have the GLM $_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma_{1.2}^2)$ with Gaussian errors and $\sigma_{1.2}^2 = \sigma_{11} - \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21}$.

Proof. Left as an exercise.

Consequently, if we have $p + q$ random variables with a joint Gaussian distribution, we may analyze p of the variables conditional upon the observed values of the other q variables. We merely use standard general linear models with the observed values of the given variables in the design matrix \mathbf{X} . The general linear model requires \mathbf{X} fixed, constant, and known without (appreciable) error. In a conditional analysis, $\mathbf{u}_{i2} = \mathbf{u}_{i02}$ and the assumption is met.

The process of *conditioning*, as used here, converts a random variable into a constant, namely the observed value. One can therefore treat a data vector either as a random vector or as a fixed constant vector by conditioning on the realized observation. In the conditional case, one makes inferences about the distribution of $\mathbf{u}_1 | \mathbf{u}_2 = \mathbf{u}_{02}$. In linear model notation, one models \mathbf{y} given \mathbf{X} at its observed value. A different \mathbf{X} matrix would produce a different conditional distribution of \mathbf{y} and possibly different results and conclusions. However, when one analyzes the unconditional joint distribution of \mathbf{u}_1 and \mathbf{u}_2 , one treats \mathbf{u}_2 as a random vector. Subsequently inferences about \mathbf{u}_1 concern its unconditional distribution, which has mean and variance different from the conditional distribution of $\mathbf{u}_1 | \mathbf{u}_2 = \mathbf{u}_{02}$.

Example 8.2 Generally speaking, analysis based on the conditional distribution restricts the scope of inferences one can draw from the data. If \mathbf{u}_1 and \mathbf{u}_2 contain heights and weights of adult humans, then the analysis of the marginal distribution of \mathbf{u}_1 would provide information about the weights of adult humans. In contrast, the analysis of the conditional distribution of $\mathbf{u}_1 | \mathbf{u}_2 = \mathbf{u}_{02}$ would provide information about the distribution of adult weights for the given set of heights.

The conditional Gaussian distribution provides the theoretical tools to support conditional analyses. If the data are approximately Gaussian in distribution, then the analysis has a sound theoretical basis. Unfortunately, the simplicity does not extend to the theory and computation of power analysis. Sampson (1974) provided an excellent discussion of the distinctions for the $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ and $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\boldsymbol{\beta}, \boldsymbol{\Sigma})$. Some discussion is also provided in later chapters on power analysis.

8.7 ASYMPTOTIC PROPERTIES

Theorem 8.9 A central limit theorem may be stated for a set of N i.i.d. length p random vectors $\{\mathbf{Y}'_i\}$, with $E(\mathbf{Y}'_i) = \boldsymbol{\mu}$ and $\mathcal{V}(\mathbf{Y}'_i) = \boldsymbol{\Sigma}$, of rank p .

(a) With $\xrightarrow{\mathcal{D}}$ indicating convergence in distribution (law), the $p \times 1$ vector

$$\bar{\mathbf{y}}_N = \sum_{i=1}^N \mathbf{Y}'_i / N = [\mathbf{Y}'_1 \ \mathbf{Y}'_2 \ \cdots \ \mathbf{Y}'_N] \mathbf{1}_N / N \tag{8.51}$$

has the property $\sqrt{N}(\bar{\mathbf{y}}_N - \boldsymbol{\mu}) \xrightarrow{\mathcal{D}} \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$ as $N \rightarrow \infty$.

(b) From the convergence in distribution, we infer an approximation for large N , namely $\sqrt{N}(\bar{\mathbf{y}}_N - \boldsymbol{\mu}) \sim \mathcal{N}_p(\mathbf{0}, \boldsymbol{\Sigma})$.

Proof. Left as an exercise.

Theorem 8.10 For $(q \times 1)$ vector $\mathbf{z}_{*N} = [f_1(\mathbf{y}_{*N})' \cdots f_q(\mathbf{y}'_{*N})]'$ = $\mathbf{f}(\mathbf{y}_{*N})$ a vector of real-valued functions of \mathbf{y}_{*N} differentiable at $\mathbf{y}_{*N} = \boldsymbol{\mu}$,

$$D = \frac{\partial \mathbf{f}'}{\partial \mathbf{y}_{*N}} \Big|_{\mathbf{y}_{*N} = \boldsymbol{\mu}} \quad (8.52)$$

exists. Equivalently, $D = \{d_{ij}\}$ and $d_{ij} = \partial f_j(\mathbf{y}_{*N}) / \partial y_{i*N}$, evaluated at $\mathbf{y}_{*N} = \boldsymbol{\mu}$. If the $(p \times 1)$ random vector $\mathbf{y}_N \stackrel{\mathcal{D}}{\rightarrow} \mathcal{N}_p(\boldsymbol{\mu}, \Sigma/N)$ has $\text{rank}(\Sigma) = p$, then random vector $\mathbf{z}_N = \mathbf{f}(\mathbf{y}_N)$ is such

$$\mathbf{z}_N \stackrel{\mathcal{D}}{\rightarrow} \mathcal{N}_q[\mathbf{f}(\boldsymbol{\mu}), D' \Sigma D / N]. \quad (8.53)$$

Proof. Mardia, Kent, and Bibby (1979) provided a proof.

8.8 THE MATRIX GAUSSIAN DISTRIBUTION

Definition 8.3 An $n \times p$ random matrix \mathbf{Y} will be said to follow a *general matrix Gaussian* distribution if and only if $\text{vec}(\mathbf{Y}) \sim (\mathcal{S})\mathcal{N}_{np}(\boldsymbol{\mu}_v, \Sigma_v)$. Necessarily $\boldsymbol{\mu}_v = \text{vec}[E(\mathbf{Y})]$ and $\mathcal{V}[\text{vec}(\mathbf{Y})] = \Sigma_v = \Sigma'_v$ is n.n.d.

Definition 8.4 The $n \times p$ random matrix \mathbf{Y} follows a *direct-product matrix Gaussian* distribution, typically abbreviated *matrix Gaussian* and written $\mathbf{Y} \sim \mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$, if and only if

- (a) $\text{vec}(\mathbf{Y}) \sim (\mathcal{S})\mathcal{N}_{np}[\text{vec}(\mathbf{M}), \Sigma \otimes \Xi]$; if and only if
- (b) $\text{vec}(\mathbf{Y}') \sim (\mathcal{S})\mathcal{N}_{n,p}[\text{vec}(\mathbf{M}'), \Xi \otimes \Sigma]$; if and only if
- (c) $\mathbf{Y} = \boldsymbol{\Psi} \mathbf{Z} \boldsymbol{\Phi}' + \mathbf{M}$ with $\text{vec}(\mathbf{Z}) \sim \mathcal{N}_{n_1 p_1}(\mathbf{0}, \mathbf{I})$ and
 - $\boldsymbol{\Psi}$ ($n \times n_1$) of rank $n_1 \geq 1$, $\Xi = \boldsymbol{\Psi} \boldsymbol{\Psi}'$,
 - $\boldsymbol{\Phi}'$ ($p_1 \times p$) of rank $p_1 \geq 1$, $\Sigma = \boldsymbol{\Phi} \boldsymbol{\Phi}'$.

Writing $\mathbf{Y} \sim \mathcal{S}\mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$ indicates $n_1 = \text{rank}(\boldsymbol{\Psi}) = \text{rank}(\Xi) < n$, or $p_1 = \text{rank}(\boldsymbol{\Phi}) = \text{rank}(\Sigma) < p$, or both.

Writing $\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$ emphasizes allowing any combination of $n_1 \leq n$ and $p_1 \leq p$.

Certain direct-product properties help in understanding the definition. As for any direct-product matrix, $\text{rank}(\Xi \otimes \Sigma) = \text{rank}(\Xi) \cdot \text{rank}(\Sigma)$. Hence both Ξ and Σ must be nonsingular for $\Xi \otimes \Sigma$ to be nonsingular (and the distribution to have a density). Furthermore the eigenvalues of $\Xi \otimes \Sigma$ are all products of the eigenvalues of Ξ and Σ of the form $\lambda_{\Xi,j} \lambda_{\Sigma,k}$. Theorem 1.5 gives $\text{vec}(\boldsymbol{\Psi} \mathbf{Z} \boldsymbol{\Phi}') = (\boldsymbol{\Phi} \otimes \boldsymbol{\Psi}) \text{vec}(\mathbf{Z})$.

The direct-product matrix Gaussian distribution has not been fully identified because, for any finite constant $a > 0$, one may write $\mathbf{Y} \sim \mathcal{N}_{n,p}[\mathbf{M}, (1/a)\Xi, a\Sigma]$.

The indeterminacy is not important for most applications because Ξ will be a known matrix and often I .

The direct-product matrix Gaussian arises naturally in many places in linear models with Gaussian errors, while the general matrix Gaussian apparently never arises naturally. Therefore, despite the ambiguity, we typically use the abbreviation “matrix Gaussian” to indicate a direct-product matrix Gaussian. The description agrees with Arnold's (1981, p. 310) discussion of the direct product matrix Gaussian as a “matrix normal” and with Gupta and Nagar's (2000), although the latter authors use slightly different notation. In a closely related approach, Mardia, Kent, and Bibby's (1979, p. 64) definition of a “normal data matrix” corresponds to the doubly special case of a direct-product matrix Gaussian with independent rows and homogeneity of mean across rows, namely $Y \sim \mathcal{N}_{n,p}(\mathbf{1}_n \boldsymbol{\mu}', I_n, \Sigma)$ for $\boldsymbol{\mu}$ a $p \times 1$ vector.

Example 8.3 A counterexample demonstrates that not all matrices of jointly Gaussian variables are direct-product matrix Gaussian. If Y is 2×2 with $\text{vec}(Y) \sim \mathcal{N}_4[0, \text{Dg}(\{1, 2, 3, 4\})]$, then Y being a direct-product matrix Gaussian requires $\{a_1, a_2, b_1, b_2\}$ exists such that $\text{Dg}(\{a_1, a_2\}) \otimes \text{Dg}(b_1, b_2) = \text{Dg}(\{1, 2, 3, 4\})$. If so, then $\{a_1 b_1 = 1, a_1 b_2 = 2, a_2 b_1 = 3, a_2 b_2 = 4\}$ implies $a_1 = 1/b_1$ and $b_2/b_1 = 2$. However, $(a_2 b_2)/(a_2 b_1) = b_2/b_1 = 4/3 \neq 2$, and the equations are inconsistent (they do not have any solution).

Example 8.4 The (direct-product) matrix Gaussian permeates the $\text{GLM}_{N,p,q}(Y_i; X_i B, \Sigma)$ with Gaussian errors. The model equation $Y = X B + E$ has $E \sim \mathcal{N}_{N,p}(\mathbf{0}, I_N, \Sigma)$ and $Y \sim \mathcal{N}_{N,p}(X B, I_N, \Sigma)$. Except for the special cases when $B = \mathbf{0}$ or $X = \mathbf{1}_N$, Y does not meet Mardia, Kent, and Bibby's definition of a normal data matrix, although E always does.

The three parameters of a (direct-product) matrix Gaussian have simple interpretations. Obviously $E(Y) = M$ is the matrix of expected values. The matrix Σ describes the covariance structure of the columns within a row. The observation may be stated precisely as follows. If d_{in} indicates an $n \times 1$ vector with 1 in position i and 0 everywhere else, then $\text{row}'_i(Y) = Y' d_{in} = \text{vec}(d'_{in} Y) = (I_p \otimes d'_{in}) \text{vec}(Y)$. The reproductive property under a linear transformation of a vector Gaussian allows writing

$$\begin{aligned} \text{row}'_i(Y) &\sim \mathcal{N}_p[(I_p \otimes d'_{in}) \text{vec}(M), (I_p \otimes d'_{in})(\Sigma \otimes \Xi)(I_p \otimes d'_{in})'] \\ &\sim \mathcal{N}_p[\text{row}'_i(M), (I_p \Sigma I_p) \otimes (d'_{in} \Xi d_{in})] \\ &\sim \mathcal{N}_p[\text{row}'_i(M), \Sigma \otimes \xi_{ii}] \\ &\sim \mathcal{N}_p[\text{row}'_i(M), \Sigma \xi_{ii}]. \end{aligned} \tag{8.54}$$

Similarly, the matrix describes the covariance structure of the rows within a column. A particular column may be written as $\text{col}_j(Y) = (I_n \otimes d'_{jp}) \text{vec}(Y')$ and

$$\begin{aligned}
\text{col}_j(\mathbf{Y}) &\sim \mathcal{N}_n \left[(\mathbf{I}_n \otimes \mathbf{d}'_{jp}) \text{vec}(\mathbf{M}'), (\mathbf{I}_n \otimes \mathbf{d}'_{jp}) (\mathbf{\Xi} \otimes \mathbf{\Sigma}) (\mathbf{I}_n \otimes \mathbf{d}'_{jp})' \right] \\
&\sim \mathcal{N}_n \left[\text{col}_j(\mathbf{M}), (\mathbf{I}_n \mathbf{\Xi} \mathbf{I}_n) \otimes (\mathbf{d}'_{jp} \mathbf{\Sigma} \mathbf{d}_{jp}) \right] \\
&\sim \mathcal{N}_n \left[\text{col}_j(\mathbf{M}), \mathbf{\Xi} \otimes \sigma_{jj} \right] \\
&\sim \mathcal{N}_n \left[\text{col}_j(\mathbf{M}), \mathbf{\Xi} \sigma_{jj} \right].
\end{aligned} \tag{8.55}$$

A useful exercise would be to define $\text{col}_j(\mathbf{Y}) = \mathbf{Y} \mathbf{d}_{jp} = (\mathbf{d}'_{jp} \otimes \mathbf{I}_n) \text{vec}(\mathbf{Y})$ and derive the distribution of the row and column in terms of $\text{vec}(\mathbf{Y})$.

Theorem 8.11 (a) With $\mathbf{T} = [t_1 \cdots t_{p_1}]$ an arbitrary real $n_1 \times p_1$ matrix, the characteristic function of $n_1 \times p_1$ $\mathbf{Z} \sim \mathcal{N}_{n_1, p_1}(\mathbf{0}, \mathbf{I}_{n_1}, \mathbf{I}_{p_1}) = \{z_{jk}\} = [z_1 \cdots z_{p_1}]$ with i.i.d. $z_{jk} \sim \mathcal{N}(0, 1)$ is $\phi(\mathbf{T}; \mathbf{Z}) = \text{E}\{\exp[\text{tr}(i\mathbf{T}'\mathbf{Z})]\} = \exp[-\text{tr}(\mathbf{T}'\mathbf{T})/2]$.

(b) For conforming constants $\mathbf{\Xi} = \mathbf{\Psi}\mathbf{\Psi}'$, $\mathbf{\Sigma} = \mathbf{\Phi}\mathbf{\Phi}'$ and \mathbf{M} , the (direct-product) matrix Gaussian $\mathbf{Y} = \mathbf{\Psi}\mathbf{Z}\mathbf{\Phi}' + \mathbf{M} \sim (S)\mathcal{N}_{n,p}(\mathbf{M}, \mathbf{\Xi}, \mathbf{\Sigma})$ has characteristic function $\phi_{\mathbf{Y}}(\mathbf{T}) = \exp[i \text{tr}(\mathbf{T}'\mathbf{M})] \exp[-\text{tr}(\mathbf{T}'\mathbf{\Xi}\mathbf{T}\mathbf{\Sigma})/2]$.

Proof. (a) Independence allows the following exchanges of operations:
 $\text{E}\{\exp[\text{tr}(i\mathbf{T}'\mathbf{Z})]\} = \text{E}\{\exp[\sum_{j=1}^{p_1} (it'_j z_j)]\} = \text{E}\{\exp[\sum_{j=1}^{p_1} \sum_{k=1}^{n_1} (it_{kj} z_{kj})]\} =$
 $\text{E}\left[\prod_{j=1}^{p_1} \prod_{k=1}^{n_1} \exp(it_{kj} z_{kj})\right] = \prod_{j=1}^{p_1} \prod_{k=1}^{n_1} \text{E}\{\exp(it_{kj} z_{kj})\} =$
 $\prod_{j=1}^{p_1} \prod_{k=1}^{n_1} \exp(-t_{kj}^2/2) = \exp\left[\sum_{j=1}^{p_1} \sum_{k=1}^{n_1} (-t_{kj}^2/2)\right] = \exp[-\text{tr}(\mathbf{T}'\mathbf{T})/2].$

(b) Lemma 7.5 gives

$$\begin{aligned}
\phi_{\mathbf{Y}}(\mathbf{T}) &= \exp[i \text{tr}(\mathbf{T}'\mathbf{M})] \exp\left\{-\text{tr}\left[\frac{(\mathbf{\Psi}'\mathbf{T}\mathbf{\Phi})'(\mathbf{\Psi}'\mathbf{T}\mathbf{\Phi})}{2}\right]\right\} \\
&= \exp[i \text{tr}(\mathbf{T}'\mathbf{M})] \exp[-\text{tr}(\mathbf{\Phi}'\mathbf{T}'\mathbf{\Psi}\mathbf{\Psi}'\mathbf{T}\mathbf{\Phi})/2] \\
&= \exp[i \text{tr}(\mathbf{T}'\mathbf{M})] \exp[-\text{tr}(\mathbf{T}'\mathbf{\Psi}\mathbf{\Psi}'\mathbf{T}\mathbf{\Phi}\mathbf{\Phi}')/2] \\
&= \exp[i \text{tr}(\mathbf{T}'\mathbf{M})] \exp[-\text{tr}(\mathbf{T}'\mathbf{\Xi}\mathbf{T}\mathbf{\Sigma})/2].
\end{aligned} \tag{8.56}$$

□

It would be hard to overstate the convenience and power of matrix Gaussian notation and describing properties of linear models. The following theorem provides one major contribution because it allows quickly deriving the distributions of estimates of expected values. Equally importantly, matrix Gaussian properties allow precisely identifying the distributions of quadratic forms and covariance matrix estimates, which lie at the heart of test statistic theory.

Theorem 8.12 If $\mathbf{Y} \sim (S)\mathcal{N}_{n,p}(\mathbf{M}, \mathbf{\Xi}, \mathbf{\Sigma})$ while $\mathbf{A} \neq \mathbf{0}$ ($n_1 \times n$), $\mathbf{B} \neq \mathbf{0}$ ($p \times p_1$), and \mathbf{C} ($n_1 \times p_1$) are finite constant matrices, then

$$\mathbf{A}\mathbf{Y}\mathbf{B} + \mathbf{C} \sim (S)\mathcal{N}_{n_1, p_1}(\mathbf{A}\mathbf{M}\mathbf{B} + \mathbf{C}, \mathbf{A}\mathbf{\Xi}\mathbf{A}', \mathbf{B}'\mathbf{\Sigma}\mathbf{B}). \tag{8.57}$$

Proof. Lemma 7.5 provides the characteristic function of a linear transformation of a random matrix. Examining the result verifies the reproductive property. \square

The theorem holds for Ξ and Σ of any rank. Here A transforms the rows and B transforms the columns. Although the CDF is always well defined, the density exists only if $\text{rank}(A\Xi A') = n_1$ and $\text{rank}(B'\Sigma B) = p_1$. However, every singular form can be expressed in terms of an embedded nonsingular one. A particularly convenient form for the singular case arises from the spectral decomposition of the covariance matrices, as in the next lemma.

Many useful results arise as special cases of the theorem. Choosing $p = p_1 = 1$ provides a standard result about the vector Gaussian. Choosing $C = 0$ and one or more other matrices as an identity matrix also produces useful special cases. The theorem helps prove the following lemma, which plays a key role in developing properties of multivariate quadratic forms.

Lemma 8.4 If $Y \sim (S)\mathcal{N}_{n,p}(M_Y, \Xi, \Sigma)$ with $\Xi = \Psi\Psi'$ of rank $n_1 \leq n$, $\Sigma = \Phi\Phi'$ of rank $p_1 \leq p$, and $M_Z = (\Psi'\Psi)^{-1}\Psi'M_Y\Phi(\Phi'\Phi)^{-1} = \Psi^+M_Y\Phi^{+t}$, then (without loss of generality) $Y = \Psi(Z + M_Z)\Phi'$ for $Z \sim \mathcal{N}_{n_1,p_1}(0, I_{n_1}, I_{p_1})$. Spectral decomposition gives $\Phi = \Upsilon\text{Dg}(\lambda)^{1/2}$ and $\Phi^{+t} = \Upsilon\text{Dg}(\lambda)^{-1/2}$.

Proof. Requiring $E(Y) = M_Y = \Psi M_Z \Phi'$ gives $(\Psi'\Psi)^{-1}\Psi'M_Y\Phi(\Phi'\Phi)^{-1} = M_Z$. Theorem 8.12 ensures Y , as a linear transformation of Z , must be matrix Gaussian with the required distribution. \square

Example 8.5 The $\text{GLM}_{N,p,q}(Y_i; X_i, B, \Sigma)$ with Gaussian errors and full-rank X has

$$\hat{B} = \left[(X'X)^{-1} X' \right] Y \sim \mathcal{N}_{q,p} \left[B, (X'X)^{-1}, \Sigma \right]. \tag{8.58}$$

If $p = 1$, then $\hat{\beta} \sim \mathcal{N}_{q,1}[\beta, (X'X)^{-1}, \sigma^2]$ if and only if $\text{vec}(\hat{\beta}) = \hat{\beta} \sim \mathcal{N}_q[\beta, (X'X)^{-1} \otimes \sigma^2]$, with $(X'X)^{-1} \otimes \sigma^2 = (X'X)^{-1} \sigma^2$. With or without full-rank X but with the requirements of full rank of $M = C(X'X)^{-1}C'$ and $C = C(X'X)^-(X'X)$ (which ensures Θ is testable),

$$\hat{\Theta} - \Theta_0 = C\hat{B}U - \Theta_0 \sim \mathcal{N}_{a,b} \left\{ \Theta - \Theta_0, [C(X'X)^{-1}C']^{-1}, U'\Sigma U \right\}. \tag{8.59}$$

If $H = X(X'X)^-X'$ and $\text{rank}(X) = r \leq q$, then $\text{rank}(H) = r$ and

$$\hat{Y} = HY \sim S\mathcal{N}_{N,p}(XB, H, \Sigma). \tag{8.60}$$

Also

$$\widehat{\mathbf{E}} = (\mathbf{I} - \mathbf{H})\mathbf{Y} \sim \mathcal{SN}_{N,p}[\mathbf{0}, (\mathbf{I} - \mathbf{H}), \boldsymbol{\Sigma}], \quad (8.61)$$

with $\text{rank}(\mathbf{I} - \mathbf{H}) = N - r$. In contrast, $\mathbf{E} \sim \mathcal{N}_{N,p}(\mathbf{0}, \mathbf{I}, \boldsymbol{\Sigma})$.

Theorem 8.13 The random matrix \mathbf{Y} is $N \times p$ with $\text{row}_i(\mathbf{Y}) = \mathbf{Y}_i$. If $(\mathbf{Y}_i)' \sim \mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is independent of $(\mathbf{Y}_{i'})'$ for $i \neq i'$ (i.i.d. rows) and $\text{rank}(\boldsymbol{\Sigma}) = p_1 \leq p$, then $\mathbf{Y} \sim \mathcal{N}_{N,p}(\mathbf{1}_N \boldsymbol{\mu}', \mathbf{I}_N, \boldsymbol{\Sigma})$. In turn, for the special case,
(a) $\text{vec}(\mathbf{Y}') = \mathbf{y} \sim \mathcal{N}_{Np}(\mathbf{1}_N \otimes \boldsymbol{\mu}, \mathbf{I}_N \otimes \boldsymbol{\Sigma})$ with $\text{rank}(\mathbf{I}_N \otimes \boldsymbol{\Sigma}) = Np_1 \leq Np$,
(b) $\text{vec}(\mathbf{Y}) = \mathbf{y} \sim \mathcal{N}_{Np}(\boldsymbol{\mu} \otimes \mathbf{1}_N, \boldsymbol{\Sigma} \otimes \mathbf{I}_N)$, and
(c) if $p_1 = p$, the density exists and may be computed as

$$\begin{aligned} f_{\mathbf{y}}(\mathbf{y}_*; \boldsymbol{\mu}, \boldsymbol{\Sigma}) &= (2\pi)^{-Np/2} |\boldsymbol{\Sigma}|^{-N/2} \exp \left[-\sum_{i=1}^N (\mathbf{Y}_{*i}' - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_{*i}' - \boldsymbol{\mu}) / 2 \right] \\ &= |2\pi \boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\text{tr} [(\mathbf{Y}_* - \mathbf{1}_N \boldsymbol{\mu}') \boldsymbol{\Sigma}^{-1} (\mathbf{Y}_* - \mathbf{1}_N \boldsymbol{\mu}')'] / 2 \right\} \\ &= |2\pi \boldsymbol{\Sigma}|^{-N/2} \exp \left\{ -\text{tr} [\boldsymbol{\Sigma}^{-1} (\mathbf{Y}_* - \mathbf{1}_N \boldsymbol{\mu}')' (\mathbf{Y}_* - \mathbf{1}_N \boldsymbol{\mu}')'] / 2 \right\}. \end{aligned} \quad (8.62)$$

Proof. Left as an exercise.

Theorem 8.14 If $\mathbf{Y} \sim \mathcal{N}_{n,p}(\mathbf{M}, \boldsymbol{\Xi}, \boldsymbol{\Sigma})$, $\text{rank}(\boldsymbol{\Xi}) = n$, and $\text{rank}(\boldsymbol{\Sigma}) = p$, then

$$\begin{aligned} f_{\mathbf{Y}}(\mathbf{Y}_*; \mathbf{M}, \boldsymbol{\Xi}, \boldsymbol{\Sigma}) &= \frac{\exp \left\{ -[\text{vec}(\mathbf{Y}_* - \mathbf{M})]' (\boldsymbol{\Xi} \otimes \boldsymbol{\Sigma})^{-1} \text{vec}(\mathbf{Y}_* - \mathbf{M}) / 2 \right\}}{(2\pi)^{np/2} |\boldsymbol{\Xi} \otimes \boldsymbol{\Sigma}|^{1/2}} \\ &= \frac{\exp \left\{ -[\text{vec}(\mathbf{Y}_* - \mathbf{M})]' (\boldsymbol{\Xi}^{-1} \otimes \boldsymbol{\Sigma}^{-1}) \text{vec}(\mathbf{Y}_* - \mathbf{M}) / 2 \right\}}{(2\pi)^{np/2} |\boldsymbol{\Xi}|^{1/2} |\boldsymbol{\Sigma}|^{1/2}} \\ &= \frac{\exp \left(-\left\{ \text{vec} [(\mathbf{Y}_* - \mathbf{M})'] \right\}' (\boldsymbol{\Sigma} \otimes \boldsymbol{\Xi})^{-1} \text{vec} [(\mathbf{Y}_* - \mathbf{M})'] / 2 \right)}{(2\pi)^{np/2} |\boldsymbol{\Sigma} \otimes \boldsymbol{\Xi}|^{1/2}} \\ &= \frac{\exp \left(-\left\{ \text{vec} [(\mathbf{Y}_* - \mathbf{M})'] \right\}' (\boldsymbol{\Sigma}^{-1} \otimes \boldsymbol{\Xi}^{-1}) \text{vec} [(\mathbf{Y}_* - \mathbf{M})'] / 2 \right)}{(2\pi)^{np/2} |\boldsymbol{\Sigma}|^{1/2} |\boldsymbol{\Xi}|^{1/2}}. \end{aligned} \quad (8.63)$$

Proof. Left as an exercise.

8.9 ASSESSING MULTIVARIATE GAUSSIAN DISTRIBUTION

Multivariate methods commonly assume the errors follow a Gaussian distribution. Whereas such methods have long been available for univariate data (since circa 1900), tests of the strong assumption of multivariate Gaussian distribution were not developed until relatively recently (circa 1970). Departure from the family of multivariate Gaussian distributions can occur in a great variety of ways. In contrast, departures from univariate Gaussian distribution occur in a

relatively small number of ways. Furthermore, obvious alternative univariate distributions often deserve consideration, such as the lognormal or Student's t . Unfortunately, any attempt to generalize the approach encounters a variety of kinds of multivariate t distributions (Kotz and Nadarajah, 2004) or a variety of other elliptical distributions. It is not clear which forms might provide reasonable alternatives. The uncertainty makes it difficult to choose sufficiently broad classes of alternatives to the null hypothesis of a multivariate Gaussian distribution.

Departure from Gaussian distribution can occur in many ways. Many multivariate methods depend on the sample covariance matrix being a good representation of associations among the variables. However, if the dependencies among some or all of the variables are not linear in nature, such as $x_2 = x_1^2 + e$, then the covariances (and associated correlations) can be very poor measures of association. Due to the variety of departures, a variety of detection techniques are needed. The techniques should include descriptive methods, graphical methods, and hypothesis tests.

Any overall test we might formulate will have to examine many features and may therefore have low sensitivity for some. Consequently if the departure from Gaussian distribution involves only a small subset of the variables, then the sensitivity of an overall test may be diluted by most of the variables being jointly Gaussian. On the other hand, any test for a specific feature or small set of features can be powerful for the feature but may not detect the departures in other unexamined features.

D'Agostino and Stephens (1986) surveyed goodness-of-fit techniques, in general. Thode (2002) discussed testing for normality in univariate and multivariate settings. Mecklin and Mundfrom (2005) provided a Monte Carlo comparison of the type I and type II error rates of tests of multivariate normality.

Graphical procedures, some as simple as scatter plots, should always be employed to visually examine the data. Two-dimensional and three-dimensional scatter plots can reveal outliers and other extreme values. Outliers can be misleading as to whether or not the data follow a Gaussian distribution. They can both conceal and falsely mimic departures. On the other hand, extreme values may not be outliers at all in the sense of errors in the data. Rather, they may indicate the need for a transformation to a Gaussian distribution. A typical sample of lognormal data will illustrate the point. Such data are usually highly skewed to the right. A few extreme points almost always occur. In general, it is often difficult to discriminate between outliers (i.e., unacceptable errors) and valid extreme values.

The difficulty of discriminating valid from invalid extreme values suggests modeling, transformation, and outlier detection should be undertaken simultaneously and interactively. In many cases “modeling” should be taken to mean *robust estimation of parameters*. Carroll and Rupert (1985) provided an excellent discussion of the principles for univariate analyses.

The $N \times p$ data matrix \mathbf{Y} has $\text{row}_i(\mathbf{Y}) = \mathbf{Y}_i$ independent of all other rows. In two or three dimensions ($p \leq 3$) we can plot the data to find outliers. Once found,

we have several options. In higher dimensions, algorithms become necessary for detecting outliers. The central question is “How far is \mathbf{Y}'_i ($p \times 1$) from the center of the cloud of data points?” A measure of distance is needed. However, estimation underlying the measure needs to be robust to the presence of several outliers masking each other.

Definition 8.5 *Masking* occurs when one or more outliers remain undetected due to the presence of other outliers.

Definition 8.6 The *Mahalanobis distance* is

$$d_i^2 = [\mathbf{Y}'_i - \mathbf{t}(\mathbf{Y})]'[\mathbf{C}(\mathbf{Y})]^{-1}[\mathbf{Y}'_i - \mathbf{t}(\mathbf{Y})], \quad (8.64)$$

with $\mathbf{t}(\mathbf{Y}) = \mathbf{Y}'\mathbf{1}_N/N = \bar{\mathbf{y}}$, the $p \times 1$ vector of arithmetic means, and $\mathbf{C}(\mathbf{Y}) = \widehat{\Sigma} = [\mathbf{Y}'\mathbf{Y}/N - \bar{\mathbf{y}}\bar{\mathbf{y}}']N/(N-1)$, the sample covariance matrix estimate.

The function suffers from the masking effect because a set of multiple outliers do not necessarily have large d_i values. Furthermore, $\mathbf{t}(\cdot)$ and $\mathcal{V}(\cdot)$ are not robust because a small cluster of outliers will attract $\mathbf{t}(\cdot)$ and inflate $\mathcal{V}(\cdot)$ in its direction.

Definition 8.7 The *breakdown point* is the number of outliers (given as a percent of number of independent sampling units) tolerated by a procedure. Larger values are better.

The Mahalanobis distance has a very low breakdown point because a few outliers can mask each other.

Definition 8.8 The *robust distance* based on the minimum volume ellipsoid (MVE) estimator, proposed by Rousseeuw and Van Zomeren (1990), is

$$RD_i = \{[\mathbf{Y}'_i - \mathbf{t}(\mathbf{Y})]'[\mathbf{S}(\mathbf{Y})]^{-1}[\mathbf{Y}'_i - \mathbf{t}(\mathbf{Y})]\}^{1/2}. \quad (8.65)$$

Here $\mathbf{t}(\mathbf{Y})$ is the MVE center and $\mathbf{S}(\mathbf{Y})$ is the corresponding covariance matrix, with both being high-breakdown estimators. Also $\mathbf{t}(\mathbf{Y})$ is the center of the MVE covering half of the observations, and $\mathbf{S}(\mathbf{Y})$ is determined by the MVE. It is multiplied by a correction factor to obtain consistency for multivariate Gaussian distributions.

8.10 TESTS FOR GAUSSIAN DISTRIBUTION

The choice of specific tests should target (1) departures anticipated as most likely and (2) departures most detrimental to the particular analysis method being used. It is advantageous if the tests suggest either a transformation which makes

the data Gaussian or an alternative data analysis method. We need tests for examining subsets because deviations from Gaussian distribution can be limited to a subset of the variables. It is difficult to check all possible subsets. Choosing subsets most likely to exhibit deviations makes it difficult to determine the true significance level of a test. Essentially three kinds of tests exist: (1) tests for marginal Gaussian distribution, (2) univariate tests of joint Gaussian distribution, and (3) multivariate tests of joint Gaussian distribution.

The simplest tests of marginal Gaussian distribution are achieved by individually testing p subhypotheses $H_{0j} : \{\mathbf{y}_j \sim \text{Gaussian}\}$ using standard testing procedures. For the overall test of $H_0 : \{\mathbf{y}_j \sim \text{Gaussian} \forall j\}$ the null hypothesis is rejected if any one of the subhypotheses is rejected. The significance level of the overall test must be controlled using principles of simultaneous testing. The method is effective in detecting the least Gaussian marginal distribution. However, it may not be powerful for detecting a subtle departure common to many or all of the variables.

Tests of marginal Gaussian distribution designed to detect subtle, common departures have been formulated in terms of measures of skewness and kurtosis. Small (1980) formulated a test in terms of skewness and kurtosis vectors, β_1 ($p \times 1$) and β_2 ($p \times 1$). Individual elements are the skewness and kurtosis parameters of the marginal distributions. Small derived scalar test statistics, Q_1 and Q_2 , each distributed as a multiple of a chi square and “nearly independent.”

One approach to detecting departures from joint Gaussian distribution is to examine the Mahalanobis distances. Clustering of observation vectors too far from (or too near to) the sample mean vector is evidence of departure from the assumption of a joint Gaussian distribution. Gnanadesikan and Kettenring (1972) discussed plotting the order statistics $\{d_{(i)}\}$ against their expected values under the null hypothesis of Gaussian distribution. The distribution of the order statistics follows a beta distribution. Small (1988) noted that the order statistics can be converted to normal scores, which should follow a Gaussian distribution under the null hypothesis. A test of univariate Gaussian distribution of the normal scores is, however, equally influenced by clusters of points too close or too far from the sample mean vector.

Another approach to accessing the joint Gaussian distribution derives from the “linear functional” characterization of Gaussian distributions. It arises from the observation that every linear combination of the variates must have a univariate Gaussian distribution.

In principle, we might test every possible linear combination for the univariate Gaussian distribution with the goal of finding the linear combination which maximizes the deviation from the Gaussian distribution. Malkovich and Afifi (1973) investigated the approach using each of three criteria: maximum skewness, maximum kurtosis, and minimum Shapiro-Wilks W statistic. If the number of variables, p , is large, then the approach will be computationally intensive. The alternative is to test specific linear combinations. If the specific linear

combinations are suggested by the data, then it is prohibitively difficult to determine the significance levels of the tests.

Cox and Small (1978) suggested yet another approach. An arbitrary transformation to a bivariate distribution, with \mathbf{a}_1 and \mathbf{a}_2 $p \times 1$, may be written as

$$\mathbf{z}_i = \begin{bmatrix} z_{1i} \\ z_{2i} \end{bmatrix} = \begin{bmatrix} \mathbf{a}'_1 \\ \mathbf{a}'_2 \end{bmatrix} \mathbf{Y}'_i = \mathbf{A} \mathbf{Y}'_i. \quad (8.66)$$

Under the Gaussian assumption, \mathbf{z}_i follows a bivariate Gaussian distribution and the conditional mean is linear in the parameters (rather than quadratic, exponential, etc.). The regression function of interest is

$$E(z_{1i}|z_{2i}) = \beta_0 + \beta_1 z_{2i} + \beta_2 z_{2i}^\lambda. \quad (8.67)$$

A test for Gaussian distribution could be obtained by testing $H_0 : \beta_2 = 0$ versus $H_A : \beta_2 \neq 0$. If $\lambda = 2$ and $\eta^2(\mathbf{a}_1, \mathbf{a}_2)$ is the sum of squares accounted for by the quadratic term in the regression model, we can maximize analytically over \mathbf{a}_1 to yield $\eta^2(\mathbf{a}_2)$, then maximize numerically over \mathbf{a}_2 to yield $\hat{\eta}_{\max}^2$. Simulations by Cox and Small (1978) allows concluding that, for $N \geq 50$, $p \leq 6$, and H_0 true, the following holds approximately:

$$\log(\hat{\eta}_{\max}^2) \sim \mathcal{N}_1[\log(5p^2/(8N)), \log^2(0.43 + 3.87/p)]. \quad (8.68)$$

Mardia (1970) defined scalar parameters for multivariate skewness and kurtosis, $\beta_{1,p}$ and $\beta_{2,p}$, proposed tests based on estimators of the parameters, and provided tables of critical values for $H_0 : \{(\beta_{1,p} = 0) \wedge [\beta_{2,p} = p(p+2)]\}$. The two subhypotheses are also of interest. Mardia (1975) proved "broadly speaking, in the presence of nonnormality, normal theory tests on *means* are influenced by $\beta_{1,p}$ whereas tests about covariances are influenced by $\beta_{2,p}$."

Definition 8.9 For a set of N vectors, each $p \times 1$, \mathbf{Y}'_i i.i.d. with mean vector $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ of rank p , the population parameters for *skewness* and *kurtosis* are, respectively,

$$\beta_{1,p} = E[(\mathbf{Y}'_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_i - \boldsymbol{\mu})]^3 \quad (8.69)$$

$$\beta_{2,p} = E[(\mathbf{Y}'_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_i - \boldsymbol{\mu})]^2.$$

The definitions are attributed to Mardia (1970).

Skewness $\beta_{1,p}$ is the expected value of the cube of the angle between vectors \mathbf{Y}'_i and \mathbf{Y}'_i (weighted by both distances) in the Mahalanobis space with metric $\boldsymbol{\Sigma}$. Kurtosis $\beta_{2,p}$ is the expected value of the squared Mahalanobis distance between vectors \mathbf{Y}'_i and \mathbf{Y}'_i in the Mahalanobis space with metric $\boldsymbol{\Sigma}$.

Theorem 8.15 For N $\mathbf{Y}'_i \sim$ i.i.d. $\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\text{rank}(\boldsymbol{\Sigma}) = p$, the population parameters for skewness and kurtosis are, respectively,

$$\beta_{1,p} = E[(\mathbf{Y}'_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_i - \boldsymbol{\mu})]^3 = E(d_{iii}^3) = 0 \tag{8.70}$$

$$\beta_{2,p} = E[(\mathbf{Y}'_i - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{Y}'_i - \boldsymbol{\mu})]^2 = E(d_{iii}^2) = p(2 + p). \tag{8.71}$$

Proof. It can be proven that the distribution of d_{iii} is symmetric about zero. Symmetry with $E(d_{iii}) = 0$ implies $E(d_{iii}^3) = 0$. Also, $E(d_{ii}^2) = p(p + 2)$ since $d_{ii}^2 \sim \chi^2(p)$. \square

Corollary 8.15.1 The scalar (univariate) Gaussian is a special case of the vector Gaussian with $p = 1$. Observation i for variable j , namely $y_{ij} \sim \mathcal{N}_1(\mu_j, \sigma_{jj})$, is i.i.d. for $i \in \{1, 2, \dots, N\}$. The corresponding population parameters for skewness and kurtosis are, respectively,

$$\beta_1 = E[(y_{ij} - \mu_j) \sigma_{jj}^{-1} (y_{ij} - \mu_j)]^3 = E(d_{iii}^3) = 0 \tag{8.72}$$

$$\beta_2 = E[(y_{ij} - \mu_j)^2 \sigma_{jj}^{-1}]^2 = E(d_{ii}^2) = 3. \tag{8.73}$$

Corollary 8.15.2 If $\mathbf{Y}'_i \sim$ i.i.d. $\mathcal{N}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\text{rank}(\boldsymbol{\Sigma}) = p$, $i \in \{1, 2, \dots, N\}$, the vectors of population parameters for skewness and kurtosis are, respectively, $\boldsymbol{\beta}_1 = \mathbf{0}$ ($p \times 1$) and $\boldsymbol{\beta}_2 = 3 \cdot \mathbf{1}$ ($p \times 1$).

The result is relevant to Small's test for marginal Gaussian distribution discussed above. Andrews, Gnanadesikan, and Warner (1971) discussed a test for multivariate Gaussian distribution in the context of transforming the data to marginal Gaussian distribution. They focused on the class of transformations proposed by Box and Cox (1964), namely $y^{(\lambda)} = (y^\lambda - 1)/\lambda$ for $\lambda \neq 0$ and $y^{(\lambda)} = \log(y)$ for $\lambda = 0$. If $\lambda = 1$, no transformation is needed. Marginal transformations to marginal Gaussian distribution could be obtained as

$$\mathbf{Y}'_i^{(\boldsymbol{\lambda})} = \left[y_{i1}^{(\lambda_1)} \ y_{i2}^{(\lambda_2)} \ \dots \ y_{ip}^{(\lambda_p)} \right]. \tag{8.74}$$

If the goal is to make every marginal distribution Gaussian, then each λ_i would be estimated separately.

Alternatively, a simultaneous estimation of $\boldsymbol{\lambda}$ can yield a transformation to joint Gaussian distribution. If $\mathbf{Y}'_i^{(\boldsymbol{\lambda})}$ is vector Gaussian for some value of $\boldsymbol{\lambda}$, then maximum likelihood methods can be used to obtain $\hat{\boldsymbol{\lambda}}$, and a likelihood ratio test can be made for $H_0 : \boldsymbol{\lambda} = \mathbf{1}$ versus the general alternative. The maximized log likelihood $L_m(\boldsymbol{\lambda})$ yields test statistic

$$2 \left[L_m(\hat{\boldsymbol{\lambda}}) - L_m(\mathbf{1}) \right] \stackrel{D}{\sim} \chi^2(p). \tag{8.75}$$

The MLE method assumes that for some value of $\boldsymbol{\lambda}$ the transformed data are

jointly Gaussian. No such value may exist because the Box-Cox transformation has limited flexibility. Even if departures occur, no value of λ may significantly improve alignment with a Gaussian distribution. Therefore, failure to reject $H_0 : \lambda = 1$ does not guarantee a Gaussian distribution. However, if the null hypothesis is rejected, then the MLE $\hat{\lambda}$ indicates a useful transformation.

Andrews et al. (1971) also considered a Box-Cox transformation of a particular projection of the data to univariate dimensions. Specifically, they suggested transforming the least Gaussian projection. First, the projection must be identified. Second, λ must be estimated for the Box-Cox transformation of the corresponding univariate distribution.

Cox and Small (1978) suggested testing $H_0 : \beta_2 = 0$ versus $H_A : \beta_2 \neq 0$ for the bivariate model

$$E(y_{ij}|y_{ij'}) = \beta_0 + \beta_1 y_{ij} + \beta_2 y_{ij}^2 \quad (8.76)$$

for variables j and j' , $j \neq j'$. The test statistic $Q_{jj'}^{(p)}$ follows a Student t distribution. By varying the choice of $j \neq j'$ the set of $p(p+1)/2$ such test statistics could be summarized in several ways. Plotting the ordered $Q_{jj'}^{(p)}$ against the expected values of such order statistics provides an example.

EXERCISES

A vector of 1's will be indicated by $\mathbf{1}$. To be explicit, one would specify the dimension, such as by writing $\mathbf{1}_N$. However, the assumption of conformation for multiplication and addition suffices to determine such dimensions. In some cases, the dimensions of such vectors, which may vary even in the same equation, are omitted in this set of exercises.

8.1 Consider $y_{ij} = a_0 + b_0 t_j + c_0 t_j^2 + a_i + b_i t_j + e_{ij}$ for $j \in \{1, \dots, p\}$ and $i \in \{1, \dots, N\}$, with $[a_i \ b_i]' \sim \mathcal{N}_2(\mathbf{0}, \mathbf{D})$, $e_{ij} \sim \mathcal{N}(0, \sigma_e^2)$ i.i.d., elements of $\{t_j, a_0, b_0, c_0\}$ are constants, and $t_j = j$. Assume that the vector $[a_i \ b_i]'$ is independent of all elements of $\{e_{ij}\}$. Also assume independence between i and i' for $i \neq i'$.

8.1.1 Completely specify the distribution of $\mathbf{y}_i = [y_{i1} \ \dots \ y_{ip}]'$.

8.1.2 Express $\mathcal{V}(y_{ij})$ as a polynomial in t_j .

8.2 Consider $y_{ij} = \mu + a_i + e_{ij}$ for $j \in \{1, 2, \dots, m\}$ and $i \in \{1, 2, \dots, N\}$, i.i.d. $a_i \sim \mathcal{N}(0, \sigma_a^2)$ independent of i.i.d. $e_{ij} \sim \mathcal{N}(0, \sigma_w^2)$ (for all i and j), $\sigma^2 = \sigma_a^2 + \sigma_w^2$, and $\rho = \sigma_a^2 / (\sigma_a^2 + \sigma_w^2)$. For $n \times m$ \mathbf{Y} , $\text{row}_i(\mathbf{Y}) = \mathbf{y}_i'$ and $\mathbf{y}_i = [y_{i1} \ \dots \ y_{im}]'$ is $m \times 1$. Also $\Sigma_i = \mathcal{V}(\mathbf{y}_i)$ is $m \times m$.

8.2.1 Completely specify the distribution of \mathbf{y}_i .

8.2.2 Give an interpretation for σ^2 .

8.2.3 Give an interpretation for ρ and specify the range of possible values of ρ .

8.2.4 Find $E(\hat{\mu})$ and $\mathcal{V}(\hat{\mu})$ for the "sample mean":

$$\hat{\mu} = \mathbf{1}'\mathbf{Y}\mathbf{1}/(nm) = (\mathbf{1}'_{nm}\mathbf{1}_{nm})^{-1}\mathbf{1}'_{nm}[\text{vec}(\mathbf{Y})] = (\mathbf{1}'_n\mathbf{1}_n)^{-1}\mathbf{1}'\mathbf{Y}\mathbf{1}(\mathbf{1}'_m\mathbf{1}_m)^{-1}.$$

8.3 Suppose i.i.d. $\mathbf{y}_i \sim \mathcal{N}_{m+1}(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ with $\boldsymbol{\mu}_i = \mu_i \mathbf{1}$ and $\boldsymbol{\Sigma}_i = \sigma^2[\rho \mathbf{1}\mathbf{1}' + (1 - \rho)\mathbf{I}]$ for $i \in \{1, \dots, N\}$. Partition \mathbf{y}_i into two subvectors, \mathbf{y}_{i1} ($m \times 1$) and y_{i2} (1×1): $\mathbf{y}_i = \begin{bmatrix} \mathbf{y}_{i1} \\ y_{i2} \end{bmatrix}$.

8.3.1 Completely specify the marginal distributions of \mathbf{y}_{i1} and of y_{i2} .

8.3.2 Completely specify the conditional distribution of $y_{i2} | \mathbf{y}_{i1} = \mathbf{u}_{i10}$.

8.3.3 Describe the behavior of $\mu_{2,1} = E(y_{i2} | \mathbf{y}_{i1} = \mathbf{u}_{i10})$ as a function of m , ρ , and σ^2 .

8.3.4 (Optional, noncredit) Describe the behavior of $\sigma_{2,1}^2 = \mathcal{V}(y_{i2} | \mathbf{y}_{i1} = \mathbf{u}_{i10})$ as a function of m , ρ and σ^2 .

8.4 Consider a matrix, \mathbf{Y} , $N \times p$, with each element of \mathbf{Y} marginally Gaussian, $y_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma^2)$, and all elements jointly Gaussian. Any pair of distinct observations has constant correlation ρ [any y_{ij} and $y_{i'j'}$ if (a) $i \neq i'$ (b) $j \neq j'$, or (c) $i \neq i'$ and $j \neq j'$].

8.4.1 What must the correlation be for the only remaining case (d) $i = i'$ and $j = j'$?

8.4.2 Clearly specify the distribution of $\text{vec}(\mathbf{Y})$.

8.4.3 Explain why or why not \mathbf{Y} is a direct-product matrix Gaussian. If it is, specify an appropriate choice of parameters.

8.5 Consider a matrix \mathbf{X} $N \times p$, with each element of \mathbf{X} marginally Gaussian, $x_{ij} \sim \mathcal{N}(\mu_{ij}, \sigma^2)$, and all elements jointly Gaussian. Also assume that any pair of distinct observations within a column (any x_{ij} and $x_{i'j}$ if $i \neq i'$) has constant correlation ρ and any pair of distinct observations within a row (any x_{ij} and $x_{ij'}$ if $j \neq j'$) are independent.

8.5.1 Clearly specify the distribution of $\text{vec}(\mathbf{X})$.

8.5.2 Explain why or why not \mathbf{X} is a direct-product matrix Gaussian. If it is, specify an appropriate choice of parameters.

8.6 Computing assignment, assuming access to SAS/IML, S+, MATLAB, or similar matrix language program. Assume all means are zero.

8.6.1 Use your knowledge of linear transformations as applied to Gaussian variables to sketch a simple algorithm to transform two independent Gaussian variables to Gaussian variables with a correlation coefficient of 0.5.

8.6.2 Based on exercise 8.6.1, generate a sample of size $N = 100$ from a bivariate normal distribution with mean zero, unit variances, and correlation $\rho = 0.5$. Provide a scatter plot of the data.

Hints

Hint 1. Review Section 1.1, especially the suggestions for writing matrices.

Hint 2. The following results may help in exercise 8.3. If $m \times m$ matrix $\mathbf{R} = \rho \mathbf{1}\mathbf{1}' + (1 - \rho)\mathbf{I}$, then the following all hold.

1. The eigenvalues of \mathbf{R} are $\lambda_1 = \rho m + (1 - \rho)$ and $\lambda_2 = (1 - \rho) = \lambda_3 = \dots = \lambda_m$.
2. $(m \times m) \mathbf{R}^{-1} = (\lambda_1 \lambda_2)^{-1}(-\rho \mathbf{1}\mathbf{1}' + \lambda_1 \mathbf{I})$.
3. A set of orthonormal eigenvectors, \mathbf{v}_i ($m \times 1$), for \mathbf{R} can be obtained easily by finding the orthonormal polynomial coefficients (Appendix, Section A.3).

4. Chapter 1 has a brief discussion of compound symmetry.
5. The spectral decomposition gives $\mathbf{R} = \mathbf{V}\text{Dg}(\boldsymbol{\lambda})\mathbf{V}' = \sum_{j=1}^m \lambda_j \mathbf{v}_j \mathbf{v}_j'$. The last form is the constituent matrix decomposition, with $\mathbf{v}_j \mathbf{v}_j' = \mathbf{G}_j$ a constituent matrix.
6. Other alternative sets of orthonormal eigenvectors can be used, such as components of a Helmert matrix.

CHAPTER 9

Univariate Quadratic Forms

9.1 MOTIVATION

In the univariate linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ estimates of $\boldsymbol{\beta}$ are *linear* transformations of the response vector \mathbf{y} . In contrast, estimates of second-order moments (including σ^2) are *quadratic* transformations of \mathbf{y} . More specifically, they are statistical *quadratic forms* in \mathbf{y} , as defined later in this chapter. Statistics for testing hypotheses are generally scalar-valued functions of quadratic forms in \mathbf{y} . The present chapter includes some of the more important properties of quadratic forms in Gaussian distributed vectors, with particular emphasis on results with applications in linear models. A series of results have individual importance in general linear univariate models. They combine to allow proving the “ANOVA theorem,” which provides the theoretical foundation for hypothesis testing in the Analysis-of-Variance and multiple regression.

9.2 CHI-SQUARE DISTRIBUTIONS

We begin by focusing on the family of central chi-square distributions and the superfamily of noncentral chi square. Noncentral distributions play an important role in power computations for tests of hypotheses. The noncentral chi square is the fundamental distribution.

Definition 9.1 If $\mathbf{z} \sim \mathcal{N}_\nu(\mathbf{0}, \mathbf{I}_\nu)$, then

$$x = \mathbf{z}'\mathbf{z} = \|\mathbf{z}\|^2 \tag{9.1}$$

has a (*central*) *chi-square distribution* with ν degrees of freedom, indicated $x \sim \chi^2(\nu)$, obviously with $0 \leq x < \infty$.

The definition is unusual. Rather than specifying a density function or some other characteristic of the distribution, we specify a random variable as a function of other random variables and define the distribution to be the distribution of the new random variable. The definition contains very little direct information about the characteristics of the distribution. We do not yet know the CDF, characteristic

function, whether it has a density function, or any other characteristic. The primary virtue of the type of definition is its simplicity. The definition also specifies the most important manner in which the distribution arises, namely, a chi square equals the sum of squares of i.i.d. $\mathcal{N}(0, 1)$ random variables. The characteristics of the central chi square are easily derived from the characteristics of the univariate Gaussian distribution. The following results are presented in many introductory texts.

Theorem 9.1 If $x \sim \chi^2(\nu)$, then its distribution is completely characterized by the following functions. The probability density function is ($\forall x_* > 0$)

$$f_x(x_*; \nu) = \frac{x_*^{(\nu-2)/2} e^{-x_*/2}}{2^{\nu/2} \Gamma(\nu/2)}. \quad (9.2)$$

The moment generating function is (for $|t| < 1/2$)

$$m_x(t; \nu) = (1 - 2t)^{-\nu/2}. \quad (9.3)$$

The characteristic function is ($\forall t$)

$$\phi_x(t; \nu) = (1 - 2it)^{-\nu/2}. \quad (9.4)$$

Proof. Johnson, Kotz, and Balakrishnan (1994) provided a detailed account.

Corollary 9.1.1 The distribution, density, and generating functions remain well-defined for any real $0 < \nu < \infty$.

Proof. Left as an exercise.

Corollary 9.1.2 If $x \sim \chi^2(\nu)$, then moments of the distribution are easily computed, for all real $m > -\nu/2$, as

$$E(x^m) = 2^m \frac{\Gamma(\nu/2 + m)}{\Gamma(\nu/2)}. \quad (9.5)$$

In particular, $E(x) = \nu$, $\mathcal{V}(x) = 2\nu$, and, if $\nu > 2$, $E(1/x) = 1/(\nu - 2)$.

Proof. Left as an exercise.

Definition 9.2 If $\mathbf{z} \sim \mathcal{N}_\nu(\boldsymbol{\mu}, \mathbf{I})$, then

$$x = \mathbf{z}'\mathbf{z} = \sum_{j=1}^{\nu} z_j^2 \quad (9.6)$$

has a *noncentral chi-square distribution* with ν degrees of freedom and noncentrality parameter $\omega = \boldsymbol{\mu}'\boldsymbol{\mu} \geq 0$, written $x \sim \chi^2(\nu, \omega)$. Alternately, if

$\mathbf{z} \sim N_p(\mathbf{1}_p, \sqrt{\omega/\nu} \mathbf{I}_p)$, then $\mathbf{z}'\mathbf{z} \sim \chi^2(\nu, \omega)$. If $\omega = 0$, the quadratic form reduces to a central chi square.

Warning! Some authors (Seattle, 1971; Johnson, Kotz, and Balakrishnan, 1995) use $\omega = \boldsymbol{\mu}'\boldsymbol{\mu}/2$ as the noncentrality parameter, while others (Rao, 1973) use $\omega = \boldsymbol{\mu}'\boldsymbol{\mu}$. Although either choice is valid, constants in the density, generating functions, and moments vary with the choice. When reading material on any function of univariate or multivariate noncentral quadratic forms, one must take care to determine which definition of noncentrality parameter is being used. The warning applies to all noncentral versions of chi square, F , and quadratic forms. The same warning also applies to noncentral versions of multivariate generalizations, including the Wishart random matrix and functions of it.

Theorem 9.2 If $x \sim \chi^2(\nu, \omega)$, then its distribution is completely characterized by the probability density function ($\forall x > 0$),

$$f_x(x_*; \nu, \omega) = \sum_{k=0}^{\infty} \frac{e^{-\omega/2} (\omega/2)^k}{k!} f_{\chi^2}(x_*; \nu + 2k, 0), \tag{9.7}$$

with $f_x(x_*; \nu + 2k, 0) = f_x(x_*; \nu + 2k)$, a central density. The moment generating function (for $|t| < 1/2$), is

$$\begin{aligned} m_x(t; \nu, \omega) &= \sum_{k=0}^{\infty} \frac{e^{-\omega/2} (\omega/2)^k}{k!} m(t; \nu + 2k, 0) \\ &= (1 - 2t)^{-\nu/2} \exp\left\{-\frac{\omega}{2} [1 - (1 - 2t)^{-1}]\right\} \\ &= (1 - 2t)^{-\nu/2} \exp\left(\frac{t\omega}{1 - 2t}\right), \end{aligned} \tag{9.8}$$

the characteristic function is ($\forall t \in \mathbb{R}$)

$$\begin{aligned} \phi_x(t; \nu, \omega) &= \sum_{k=0}^{\infty} \frac{e^{-\omega/2} (\omega/2)^k}{k!} \phi(t; \nu + 2k, 0) \\ &= (1 - 2it)^{-\nu/2} \exp\left\{-\frac{\omega}{2} [1 - (1 - 2it)^{-1}]\right\} \\ &= (1 - 2it)^{-\nu/2} \exp\left(\frac{it\omega}{1 - 2it}\right), \end{aligned} \tag{9.9}$$

the cumulant generating function is (for $|t| < 1/2$),

$$c_x(t; \nu, \omega) = -(\nu/2)\log(1 - 2t) + t\omega/(1 - 2t), \tag{9.10}$$

and cumulant m is

$$\kappa_m(x) = 2^{m-1} (m - 1)! (\nu + m\omega). \tag{9.11}$$

Proof. If $\mathbf{z} \sim \mathcal{N}_\nu(\boldsymbol{\mu}, \mathbf{I}_\nu)$ and $x = \mathbf{z}'\mathbf{z}$, the MGF is

$$\begin{aligned} m_x(t; \nu, \omega) &= \mathbb{E}[\exp(t\mathbf{z}'\mathbf{z})] \\ &= (2\pi)^{-\nu/2} \int_{\mathbb{R}^\nu} \exp[t\mathbf{z}'_*\mathbf{z}_* - (\mathbf{z}_* - \boldsymbol{\mu})'I(\mathbf{z}_* - \boldsymbol{\mu})/2] d\mathbf{z}_* \\ &= (2\pi)^{-\nu/2} \int_{\mathbb{R}^\nu} \exp[t\mathbf{z}'_*\mathbf{z}_* - \mathbf{z}'_*\mathbf{z}_*/2 + \boldsymbol{\mu}'\mathbf{z}_* - \boldsymbol{\mu}'\boldsymbol{\mu}/2] d\mathbf{z}_* \\ &= (2\pi)^{-\nu/2} \int_{\mathbb{R}^\nu} \exp[-(1-2t)\mathbf{z}'_*\mathbf{z}_*/2 + \boldsymbol{\mu}'\mathbf{z}_* - \boldsymbol{\mu}'\boldsymbol{\mu}/2] d\mathbf{z}_*. \quad (9.12) \end{aligned}$$

If $\mathbf{A}^{-1} = (1 - 2t)\mathbf{I}_\nu$ and $\omega = \boldsymbol{\mu}'\boldsymbol{\mu}$, then

$$\begin{aligned} m_x(t; \nu, \omega) &= (2\pi)^{-\nu/2} e^{-\omega/2} \int_{\mathbb{R}^\nu} \exp(-\mathbf{z}'_*\mathbf{A}^{-1}\mathbf{z}_*/2 + \boldsymbol{\mu}'\mathbf{z}_*) d\mathbf{z}_* \\ &= e^{-\omega/2} |\mathbf{A}|^{1/2} \int_{\mathbb{R}^\nu} \exp(\boldsymbol{\mu}'\mathbf{z}_*) \left[(2\pi)^{-\nu/2} |\mathbf{A}|^{-1/2} \exp(-\mathbf{z}'_*\mathbf{A}^{-1}\mathbf{z}_*/2) \right] d\mathbf{z}_*. \quad (9.13) \end{aligned}$$

The last integral is the MGF of a $\mathcal{N}_\nu(\mathbf{0}, \mathbf{A})$, with density contained in the brackets (and with $\boldsymbol{\mu}$ replacing the usual \boldsymbol{t}). Therefore

$$\begin{aligned} m_x(t; \nu, \omega) &= e^{-\omega/2} |\mathbf{A}|^{1/2} \exp(\boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}/2) \\ &= (1 - 2t)^{-\nu/2} \exp\{[(1 - 2t)^{-1}\boldsymbol{\mu}'\boldsymbol{\mu} - \omega]/2\} \\ &= (1 - 2t)^{-\nu/2} \exp[t\omega/(1 - 2t)]. \quad (9.14) \end{aligned}$$

The characteristic function of the distribution with the density shown is

$$\begin{aligned} \phi_x(t; \nu, \omega) &= \mathbb{E}(e^{itx}) = \int_0^\infty e^{itx} f_x(x_*; \nu, \omega) dx_* \\ &= \int_0^\infty e^{itx} \sum_{k=0}^\infty \frac{e^{-\omega/2} (\omega/2)^k}{k!} f_{x_*}(x_*; \nu + 2k, 0) dx_*. \quad (9.15) \end{aligned}$$

The bounded convergence theorem permits an interchange of the infinite summation and the integral:

$$\begin{aligned} \phi_x(t; \nu, \omega) &= \sum_{k=0}^\infty \frac{e^{-\omega/2} (\omega/2)^k}{k!} \int_0^\infty e^{itx} f_x(x_*; \nu + 2k, 0) dx_* \\ &= \sum_{k=0}^\infty \frac{e^{-\omega/2} (\omega/2)^k}{k!} \phi_x(t; \nu + 2k, 0) \\ &= \sum_{k=0}^\infty \frac{e^{-\omega/2} (\omega/2)^k}{k!} (1 - 2it)^{-(\nu+2k)/2} \\ &= (1 - 2it)^{-\nu/2} e^{-\omega/2} \sum_{k=0}^\infty \frac{1}{k!} \left(\frac{\omega/2}{1 - 2it} \right)^k. \quad (9.16) \end{aligned}$$

Recalling $e^x = \sum_{k=0}^\infty x^k/k!$ and simplifying give the final form. \square

The derivation of the characteristic function corresponding to the density can be repeated with $\phi_x(t; \nu, \omega)$ replaced by $m_x(t; \nu, \omega)$ and omitted to prove the stated PDF corresponds to the stated MGF. Since we have previously derived the MGF from basic principles, the verification proves the stated PDF is correct.

The PDFs in the expression for the density are weighted by Poisson probabilities for $k \in \{0, 1, 2, \dots\}$, which means the density describes a mixture. Similarly, the noncentral MGF and characteristic function are weighted averages of central MGFs and characteristic functions.

Corollary 9.2.1 The distribution, density, and generating functions remain well defined for any real $0 < \nu < \infty$.

Although Siegel (1979) discussed a distribution with zero degrees of freedom, we do not consider it here. Johnson, Kotz, and Balakrishnan (1994) provided related discussions.

Corollary 9.2.2 If $\mathbf{z} \sim \mathcal{N}_\nu(\mathbf{0}, \mathbf{I}_\nu)$, then $\mathbf{z}'\mathbf{z} \sim \chi^2(\nu, 0)$, and equivalently $\mathbf{z}'\mathbf{z} \sim \chi^2(\nu)$, i.e., the central chi-square is a special case of the noncentral chi-square with $\omega = 0$.

Corollary 9.2.3 If $\mathbf{y} \sim \mathcal{N}_\nu(\boldsymbol{\mu}, \sigma^2 \mathbf{I}_\nu)$, then $\mathbf{y}'\mathbf{y}/\sigma^2 \sim \chi^2(\nu, \boldsymbol{\mu}'\boldsymbol{\mu}/\sigma^2)$.

Corollary 9.2.4 For any finite set of independent chi-square random variables $\{x_j\}$, with $x_j \sim \chi^2(\nu_j, \omega_j)$,

$$\sum_{j=1}^n x_j \sim \chi^2\left(\sum_{j=1}^n \nu_j, \sum_{j=1}^n \omega_j\right). \tag{9.17}$$

Corollary 9.2.5 If $x \sim \chi^2(\nu, \omega)$, then $E(x) = \nu + \omega$ and $\mathcal{V}(x) = 2\nu + 4\omega$.

Proof. Left as exercises.

9.3 GENERAL PROPERTIES OF QUADRATIC FORMS

As briefly discussed in Chapter 1, in the study of matrix algebra a quadratic form is $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ for any conforming \mathbf{A} and \mathbf{y} . Without loss of generality, \mathbf{A} may be assumed to be symmetric (Lemma 1.4) because $q = \mathbf{y}'\mathbf{A}\mathbf{y} = \mathbf{y}'\mathbf{B}\mathbf{y}$ with $\mathbf{B} = [(\mathbf{A} + \mathbf{A}')/2]$. The result allows taking advantage of the special properties of symmetric matrices. Most importantly, we are assured the middle matrix has a spectral decomposition, $\mathbf{B} = \mathbf{V}_B \text{Dg}(\boldsymbol{\lambda}_B) \mathbf{V}'_B$, with square, full-rank, and orthonormal \mathbf{V}_B .

In studying quadratic forms, it is helpful to remember the eigenvalues of \mathbf{B} are necessarily real, but they may be positive, negative, or zero. When no negative

eigenvalues occur, the middle matrix in a quadratic form can be expressed as $\mathbf{A} = \mathbf{M}'\mathbf{M}$. Many special properties of inner and outer product matrices are mentioned in Chapter 1. Such products are automatically symmetric. They never have any negative eigenvalues, although some may be zero. In all cases $\text{rank}(\mathbf{M}'\mathbf{M}) = \text{rank}(\mathbf{M}\mathbf{M}') = \text{rank}(\mathbf{M})$.

Definition 9.3 With $n \times 1$ random vector \mathbf{y} and $n \times n$ constant middle matrix \mathbf{A} , the scalar $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ is a random *quadratic form*.

The next result is true for \mathbf{y} following any distribution with finite second moments. For comparison recall $E(\mathbf{y}\mathbf{y}') = \Sigma + \mu\mu'$.

Theorem 9.3 If \mathbf{y} ($n \times 1$) is *any* random vector with finite mean $E(\mathbf{y}) = \mu$ ($n \times 1$) and finite dispersion $\mathcal{V}(\mathbf{y}) = \Sigma$ ($n \times n$), then for any symmetric, finite constant matrix \mathbf{A} ($n \times n$)

$$E(\mathbf{y}'\mathbf{A}\mathbf{y}) = \text{tr}(\mathbf{A}\Sigma) + \mu'\mathbf{A}\mu. \quad (9.18)$$

Proof. The cyclical property of the trace for conforming matrices, $\text{tr}(\mathbf{ABC}) = \text{tr}(\mathbf{CAB}) = \text{tr}(\mathbf{BCA})$, applies. Since the quadratic form is a scalar, it is equal to its trace. Therefore $E(\mathbf{y}'\mathbf{A}\mathbf{y}) = E[\text{tr}(\mathbf{y}'\mathbf{A}\mathbf{y})] = E[\text{tr}(\mathbf{A}\mathbf{y}\mathbf{y}')] = \text{tr}[\mathbf{A}E(\mathbf{y}\mathbf{y}')] = \text{tr}[\mathbf{A}(\Sigma + \mu\mu')] = \text{tr}(\mathbf{A}\Sigma) + \mu'\mathbf{A}\mu. \quad \square$

Corollary 9.3 If $\mathbf{A} \neq \mathbf{A}'$ the result still holds.

Proof. Although Lemma 1.4 applies, $\text{tr}(\mathbf{B}\Sigma) = \text{tr}\{[(\mathbf{A} + \mathbf{A}')/2]\Sigma\} = [\text{tr}(\mathbf{A}\Sigma) + \text{tr}(\mathbf{A}'\Sigma)]/2$ gives one pause. However, $\text{tr}(\mathbf{A}\Sigma) = \text{tr}(\Sigma\mathbf{A}') = \text{tr}(\mathbf{A}'\Sigma)$ proves the result by trace invariance to transposition and permutation. \square

9.4 PROPERTIES OF QUADRATIC FORMS IN GAUSSIAN VECTORS

Definition 9.4 A (univariate) *quadratic form in Gaussian variables* is a (scalar) random variable $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ for $n \times n$ constant (finite) \mathbf{A} with $\text{rank}(\mathbf{A}\Sigma) > 0$ and $\mathbf{y} \sim (\mathcal{S})\mathcal{N}_n(\mu, \Sigma)$. Without loss of generality $\mathbf{A} = \mathbf{A}'$.

Example 9.1 Although not obvious, q may be negative. More precisely, although $\Pr\{\mathbf{y}'\mathbf{y} < 0\} = 0$, depending on the choice of \mathbf{A} , it may be that $\Pr\{\mathbf{y}'\mathbf{A}\mathbf{y} < 0\} > 0$. If, as an example, $n = 1$, $\mathbf{A} = [-1]$, and $\Sigma = 1$, then $q = -\mathbf{y}'\mathbf{y} = -y_1^2$ and $\Pr\{q < 0\} = 1$.

Example 9.2 If $\mathbf{A} = \mathbf{I}_n$ and $\mathbf{y} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n\sigma^2)$, then

$$q = \mathbf{y}'\mathbf{y} = \sigma^2 \sum_{j=1}^n y_j^2 / \sigma^2 = \sigma^2 \sum_{j=1}^n x_j. \tag{9.19}$$

Here $x_j \sim \chi^2(1, 0)$ and $x_j \perp x_{j'}$ if $j \neq j'$ due to the independence of the underlying Gaussian variables. In turn $q/\sigma^2 \sim \chi^2(n, 0)$.

Example 9.3 If $\mathbf{A} = \mathbf{I}_n$ and $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \mathbf{I}_n\sigma^2)$, then

$$q = \mathbf{y}'\mathbf{y} = \sigma^2 \sum_{j=1}^n y_j^2 / \sigma^2 = \sigma^2 \sum_{j=1}^n x_j. \tag{9.20}$$

Here $x_j \sim \chi^2(1, \mu_j^2/\sigma^2)$ and $x_j \perp x_{j'}$ if $j \neq j'$ due to the independence of the underlying Gaussian variables. In turn $q/\sigma^2 \sim \chi^2(n, \boldsymbol{\mu}'\boldsymbol{\mu}/\sigma^2)$.

The last two examples derive expressions for special quadratic forms in terms of simple sets of underlying Gaussian and chi square variables. In the most general case, similar expressions can be found in terms of weighted sums of possibly noncentral chi-square variables. The following three theorems provide explicit decompositions for increasing more general quadratic forms.

Theorem 9.4 Random $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \sigma^2\mathbf{I}_n)$ and $n \times n$ constant (finite) $\mathbf{A} = \mathbf{A}'$ of rank $0 < n_1 \leq n$ define $q = \mathbf{y}'\mathbf{A}\mathbf{y}$. Spectral decomposition gives $\mathbf{A} = \mathbf{V}_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{V}_1'$, with orthonormal-by-column $\mathbf{V}_1 = [\mathbf{v}_{1,1} \cdots \mathbf{v}_{1,n_1}]$ ($n \times n_1$) and $\boldsymbol{\lambda}_1$ ($n_1 \times 1$).

(a) In any such setting

$$q = \sigma^2 \sum_{k=1}^{n_1} \lambda_{1,k} x_k, \tag{9.21}$$

with $\{x_k\}$ mutually independent, $x_k \sim \chi^2(1, \omega_k)$, and $\omega_k = (\mathbf{v}'_{1,k}\boldsymbol{\mu})^2 / \sigma^2 = \boldsymbol{\mu}'\mathbf{v}_{1,k}\mathbf{v}'_{1,k}\boldsymbol{\mu} / \sigma^2 = \mathbf{v}'_{1,k}\boldsymbol{\mu}\boldsymbol{\mu}'\mathbf{v}_{1,k} / \sigma^2$.

(b) With the additional requirement of $\mathbf{A} = \mathbf{A}^2$ (idempotent with rank n_1), if

$$s = \sum_{k=1}^{n_1} x_k \tag{9.22}$$

$$\omega_+ = \sum_{k=1}^{n_1} \omega_k = \boldsymbol{\mu}'\mathbf{V}_1\mathbf{V}_1'\boldsymbol{\mu} / \sigma^2 \tag{9.23}$$

then $q = \sigma^2 s$ and $s \sim \chi^2(n_1, \omega_+)$.

Proof. The results are special cases of the next theorem. However, the simpler proof for the special case more clearly illustrates the principles involved. For the special case (a), without loss of generality $\mathbf{y} = \sigma(\mathbf{z} + \boldsymbol{\mu}/\sigma)$ for $\mathbf{z} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{I}_n)$.

As an orthonormal linear transformation of a unit vector Gaussian, by Lemma 8.2, $V_1'z \sim \mathcal{N}_{n_1}(\mathbf{0}, I_{n_1})$. In turn, if $\mu_1 = V_1'\mu/\sigma$, then $\mathbf{x} = (z_1 + \mu_1) \sim \mathcal{N}_{n_1}(\mu_1, I_{n_1})$. Furthermore

$$\begin{aligned} q &= \mathbf{y}'\mathbf{A}\mathbf{y} = \sigma(z + \mu/\sigma)'V_1\text{Dg}(\lambda_1)V_1'(z + \mu/\sigma)\sigma \\ &= (V_1'z + V_1'\mu/\sigma)'\text{Dg}(\sigma^2\lambda_1)(V_1'z + V_1'\mu/\sigma) \\ &= (z_1 + \mu_1)'\text{Dg}(\sigma^2\lambda_1)(z_1 + \mu_1) \\ &= \mathbf{x}'\text{Dg}(\sigma^2\lambda_1)\mathbf{x} \\ &= \sum_{k=1}^{n_1} \sigma^2\lambda_{1,k}x_k^2. \end{aligned} \tag{9.24}$$

Here $\{x_k^2\}$ are independent and $x_k^2 \sim \chi^2(1, \omega_k)$.

For the special case (b), $\mathbf{A} = \mathbf{A}^2$ implies $\lambda_{1,k} \equiv 1$. □

The following theorem contains the preceding two as special cases. It also covers the two examples discussed at the beginning of the present section.

Theorem 9.5 If $\mathbf{y} \sim \mathcal{N}_n(\mu_y, \Sigma)$ with $\text{rank}(\Sigma) = n$ and $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ for $n \times n$ constant (finite) $\mathbf{A} = \mathbf{A}'$ of rank $0 < n_1 \leq n$, $\Sigma = \Phi\Phi'$ with Φ $n \times n$ of rank n , then $\mathbf{B} = \Phi'\mathbf{A}\Phi$ is $n \times n$, symmetric, and rank n_1 . Spectral decomposition gives $\mathbf{B} = V_1\text{Dg}(\lambda_1)V_1'$, with columnwise orthonormal $V_1 = [v_{1,1} \cdots v_{1,n_1}]$ of dimension $n \times n_1$ and λ_1 $n_1 \times 1$. Furthermore

$$q = \sum_{k=1}^{n_1} \lambda_{1,k}x_k, \tag{9.25}$$

with $\{x_k\}$ mutually independent, $x_k \sim \chi^2(1, \omega_k)$, and $\omega_k = (v'_{1,k}\Phi^{-1}\mu_y)^2 = \mu_y'\Phi^{-t}v_{1,k}v'_{1,k}\Phi^{-1}\mu_y = v'_{1,k}\Phi^{-1}\mu_y\mu_y'\Phi^{-t}v_{1,k}$.

Proof. The result is a special case of the next theorem. However, the simpler proof for the special case more clearly illustrates the principles involved. For the special case, if $\mu_z = \Phi^{-1}\mu_y$, without loss of generality, $\mathbf{y} = \Phi(z + \mu_z)$ with $z \sim \mathcal{N}_n(\mathbf{0}, I_n)$. Hence

$$\begin{aligned} q &= \mathbf{y}'\mathbf{A}\mathbf{y} = [\Phi(z + \mu_z)]'\mathbf{A}[\Phi(z + \mu_z)] \\ &= (z + \mu_z)\Phi'\mathbf{A}\Phi(z + \mu_z) \\ &= (z + \mu_z)V_1\text{Dg}(\lambda_1)V_1'(z + \mu_z) \\ &= (V_1'z + V_1'\mu_z)'\text{Dg}(\lambda_1)(V_1'z + V_1'\mu_z) \\ &= \mathbf{x}'\text{Dg}(\lambda_1)\mathbf{x} \\ &= \sum_{k=1}^{n_1} \lambda_{1,k}x_k^2. \end{aligned} \tag{9.26}$$

As an orthonormal linear transformation of a unit vector Gaussian, by Lemma 8.2, $V_1'z \sim \mathcal{N}_{n_1}(\mathbf{0}, I_{n_1})$. In turn, $\mathbf{x} = (V_1'z + V_1'\mu_z) \sim \mathcal{N}_{n_1}(V_1'\mu_z, I_{n_1})$ is a set of

independent Gaussian variables with $E(x_k) = \mathbf{v}'_{1k}\boldsymbol{\mu}_z$. In turn $\{x_k^2\}$ are independent and $x_k^2 \sim \chi^2(1, \omega_k)$, with $\omega_k = (\mathbf{v}'_{1k}\boldsymbol{\mu}_z)^2 = (\mathbf{v}'_{1k}\boldsymbol{\Phi}^{-1}\boldsymbol{\mu}_y)^2$. \square

Theorem 9.6 If $\mathbf{y} \sim (\mathcal{S})\mathcal{N}_n(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\text{rank}(\boldsymbol{\Sigma}) = n_1$ for $1 \leq n_1 \leq n$ and $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ for $n \times n$ constant (finite) \mathbf{A} , then the following hold.

(a) Without loss of generality, $\mathbf{A} = \mathbf{A}'$ and $\boldsymbol{\Sigma} = \boldsymbol{\Phi}\boldsymbol{\Phi}'$, with $\boldsymbol{\Phi}$ $n \times n_1$ of rank n_1 , while $\boldsymbol{\Phi}^+ = (\boldsymbol{\Phi}'\boldsymbol{\Phi})^{-1}\boldsymbol{\Phi}'$.

(b) The matrix $\mathbf{B} = \boldsymbol{\Phi}'\mathbf{A}\boldsymbol{\Phi}$ is $n_1 \times n_1$ and symmetric with $\text{rank}(\mathbf{B}) = n_2 \leq n_1 \leq n$. If $n_2 > 0$, then $\mathbf{B} = \mathbf{V}_2\text{Dg}(\boldsymbol{\lambda}_2)\mathbf{V}_2'$, with columnwise orthonormal $\mathbf{V}_2 = [\mathbf{v}_{2,1} \cdots \mathbf{v}_{2,n_2}]$ of dimension $n_1 \times n_2$ and $\boldsymbol{\lambda}_2$ $n_2 \times 1$.

(c) Furthermore

$$q = \sum_{k=1}^{n_2} \lambda_{2,k} x_k, \quad (9.27)$$

with $\{x_k\}$ mutually independent, $x_k \sim \chi^2(1, \omega_k)$, and $\omega_k = (\mathbf{v}'_{2,k}\boldsymbol{\Phi}^+\boldsymbol{\mu}_y)^2 = \boldsymbol{\mu}'_y(\boldsymbol{\Phi}^+)' \mathbf{v}_{2,k} \mathbf{v}'_{2,k} \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y$.

Proof. With $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I}_{n_1})$, Lemma 8.3 ensures $n_1 \times 1$ constant $\boldsymbol{\mu}_z = \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y$ exists such that $\mathbf{y} = \boldsymbol{\Phi}(\mathbf{z} + \boldsymbol{\mu}_z)$. In turn

$$\begin{aligned} q &= \mathbf{y}'\mathbf{A}\mathbf{y} = (\mathbf{z} + \boldsymbol{\mu}_z)' \boldsymbol{\Phi}'\mathbf{A}\boldsymbol{\Phi}(\mathbf{z} + \boldsymbol{\mu}_z) \\ &= (\mathbf{z} + \boldsymbol{\mu}_z)' \mathbf{B}(\mathbf{z} + \boldsymbol{\mu}_z). \end{aligned} \quad (9.28)$$

Symmetry of \mathbf{A} ensures the symmetry of $\mathbf{B} = \boldsymbol{\Phi}'\mathbf{A}\boldsymbol{\Phi}$. Hence \mathbf{B} has $\text{rank}(\mathbf{B}) = n_2$ nonzero eigenvalues $\{\lambda_{2,k}\}$ (all real and any mixture of positive and negative values), with $0 < n_2 \leq n_1$. Also \mathbf{B} has $n_1 - n_2$ zero eigenvalues. The spectral decomposition allows writing

$$\mathbf{B} = \mathbf{V}_2\text{Dg}(\boldsymbol{\lambda}_2)\mathbf{V}_2'. \quad (9.29)$$

Here $\text{Dg}(\boldsymbol{\lambda}_2)$ is $n_2 \times n_2$, while \mathbf{V}_2 is the $n_1 \times n_2$ columnwise orthonormal matrix of corresponding eigenvectors. In turn

$$\begin{aligned} q &= (\mathbf{z} + \boldsymbol{\mu}_z)' \mathbf{V}_2\text{Dg}(\boldsymbol{\lambda}_2)\mathbf{V}_2'(\mathbf{z} + \boldsymbol{\mu}_z) \\ &= [\mathbf{V}_2'(\mathbf{z} + \boldsymbol{\mu}_z)]' \text{Dg}(\boldsymbol{\lambda}_2) [\mathbf{V}_2'(\mathbf{z} + \boldsymbol{\mu}_z)]. \end{aligned} \quad (9.30)$$

The vector

$$\mathbf{u} = \mathbf{V}_2'(\mathbf{z} + \boldsymbol{\mu}_z) = \mathbf{V}_2'\mathbf{z} + \mathbf{V}_2'\boldsymbol{\mu}_z \quad (9.31)$$

is $n_2 \times 1$. Lemma 8.2, with $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I}_{n_1})$, gives $\mathbf{V}_2'\mathbf{z} \sim \mathcal{N}_{n_2}(\mathbf{0}, \mathbf{I}_{n_2})$. In turn

$$\mathbf{u} = \mathbf{V}_2'\mathbf{z} + \mathbf{V}_2'\boldsymbol{\mu}_z \sim \mathcal{N}_{n_2}(\mathbf{V}_2'\boldsymbol{\mu}_z, \mathbf{I}_{n_2}). \quad (9.32)$$

Therefore $\{u_k\}$ are mutually independent with $u_k \sim \mathcal{N}_1(\mu_{2,k}, 1)$ and

$$\mu_{2,k} = \mathbf{v}'_{2,k}\boldsymbol{\mu}_z = \mathbf{v}'_{2,k}\boldsymbol{\Phi}^+\boldsymbol{\mu}_y. \quad (9.33)$$

Furthermore

$$q = \mathbf{u}'\text{Dg}(\boldsymbol{\lambda}_2)\mathbf{u} = \sum_{k=1}^{n_2} \lambda_{2,k} u_k^2, \quad (9.34)$$

with $\{u_k^2\}$ being mutually independent and $u_k^2 \sim \chi^2(1, \mu_{2,k}^2)$. \square

Corollary 9.6.1 The characteristic function, moment generating function, cumulant generating function, and cumulant m of q are

$$\begin{aligned} \phi_q(t; \boldsymbol{\lambda}_2, \boldsymbol{\omega}) &= \prod_{k=1}^{n_2} \left[(1 - 2i\lambda_{2,k}t)^{-1/2} \exp\left(\frac{it\lambda_{2,k}\omega_k}{1 - 2i\lambda_{2,k}t}\right) \right] \\ &= \left[\prod_{k=1}^{n_2} (1 - 2i\lambda_{2,k}t)^{-1/2} \right] \exp\left(\sum_{k=1}^{n_2} \frac{it\lambda_{2,k}\omega_k}{1 - 2i\lambda_{2,k}t}\right) \end{aligned} \quad (9.35)$$

$$\begin{aligned} m_q(t; \boldsymbol{\lambda}_2, \boldsymbol{\omega}) &= \prod_{k=1}^{n_2} \left[(1 - 2\lambda_{2,k}t)^{-1/2} \exp\left(\frac{t\lambda_{2,k}\omega_k}{1 - 2\lambda_{2,k}t}\right) \right] \\ &= \left[\prod_{k=1}^{n_2} (1 - 2\lambda_{2,k}t)^{-1/2} \right] \exp\left(\sum_{k=1}^{n_2} \frac{t\lambda_{2,k}\omega_k}{1 - 2\lambda_{2,k}t}\right) \end{aligned} \quad (9.36)$$

$$\begin{aligned} c_q(t; \boldsymbol{\lambda}_2, \boldsymbol{\omega}) &= \sum_{k=1}^{n_2} \left[-\frac{1}{2} \log(1 - 2\lambda_{2,k}t) + \left(\frac{t\lambda_{2,k}\omega_k}{1 - 2\lambda_{2,k}t}\right) \right] \\ &= -\frac{1}{2} \sum_{k=1}^{n_2} \log(1 - 2\lambda_{2,k}t) + \sum_{k=1}^{n_2} \frac{t\lambda_{2,k}\omega_k}{1 - 2\lambda_{2,k}t} \end{aligned} \quad (9.37)$$

$$\begin{aligned} \kappa_m(q; \boldsymbol{\lambda}_2, \boldsymbol{\omega}) &= \sum_{k=1}^{n_2} [\lambda_{2,k}^m 2^{m-1} (m-1)! (1 + m\omega_k)] \\ &= 2^{m-1} (m-1)! \sum_{k=1}^{n_2} [\lambda_{2,k}^m (1 + m\omega_k)]. \end{aligned} \quad (9.38)$$

Proof. The characteristic function of a sum of independent random variables is the product of the individual (marginal) characteristic functions. As a linear transformation of a noncentral chi square, the characteristic function of each x_k can be found by applying the form $\phi_{ay+b}(t) = \exp(itb)\phi_x(at)$ to the characteristic function of a $\chi^2(1, \omega_k)$. A parallel approach applies to the moment generating function. The cumulant generating function is available by taking a logarithm. The cumulant reflects the simple impact of a linear transformation on cumulants and the fact that the cumulant of a sum of independent random variables is the sum of the individual (marginal) cumulants. \square

Corollary 9.6.2 If $\lambda_{2,k} \equiv \lambda_1$, then $q/\lambda_1 \sim \chi^2(n_2, \sum_{k=1}^{n_2} \omega_k)$.

Proof. The characteristic function of q reduces to

$$\begin{aligned} \phi_q(t; \mathbf{1}_{n_2} \lambda_1, \boldsymbol{\omega}) &= \left[\prod_{k=1}^{n_2} (1 - 2i\lambda_1 t)^{-1/2} \right] \exp \left(\sum_{k=1}^{n_2} \frac{it\lambda_1 \omega_k}{1 - 2i\lambda_1 t} \right) \\ &= (1 - 2i\lambda_1 t)^{-n_2/2} \exp \left(\frac{it\lambda_1}{1 - 2i\lambda_1 t} \sum_{k=1}^{n_2} \omega_k \right). \end{aligned} \quad (9.39)$$

Observing the effect of transforming q to q/λ_1 completes the proof. \square

The theorem characterizes a quadratic form in Gaussian random variables regardless of whether or not the quadratic form has a chi-square distribution. As will be seen in studying estimates of variability in linear models, perhaps the most important univariate quadratic forms are scaled chi-square random variables. However, many important quadratic forms are not scaled chi squares. They occur in a variety of settings, including variance component models and the distributions of test statistics.

Corollary 9.6.3 The mean is

$$\begin{aligned} E(q) &= \sum_{k=1}^{n_2} \lambda_{2,k} + \sum_{k=1}^{n_2} \lambda_{2,k} \omega_k \\ &= \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \sum_{k=1}^{n_2} \lambda_{2,k} \boldsymbol{\mu}'_y (\boldsymbol{\Phi}^+)' \mathbf{v}_{2,k} \mathbf{v}'_{2,k} \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y \\ &= \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}'_y (\boldsymbol{\Phi} \boldsymbol{\Phi}^+)' \mathbf{A} (\boldsymbol{\Phi} \boldsymbol{\Phi}^+) \boldsymbol{\mu}_y. \end{aligned} \quad (9.40)$$

If $\boldsymbol{\Sigma}$ has full rank of n , then the last form reduces to $E(q) = \text{tr}(\mathbf{A}\boldsymbol{\Sigma}) + \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$ and coincides with the result in Theorem 9.3 (which holds whether or not \mathbf{y} is Gaussian and whether or not $\boldsymbol{\Sigma}$ is full rank). Also, with $\mathbf{B} = \boldsymbol{\Phi}' \mathbf{A} \boldsymbol{\Phi}$, $\boldsymbol{\Phi} = \mathbf{V}_2 \text{Dg}(\boldsymbol{\lambda}_2)^{1/2}$ and $\boldsymbol{\Phi}^+ = \text{Dg}(\boldsymbol{\lambda}_2)^{-1/2} \mathbf{V}'_2$, and $\boldsymbol{\Phi} \boldsymbol{\Phi}^+ = \mathbf{V}_2 \mathbf{V}'_2$. Furthermore

$$\begin{aligned} \mathcal{V}(q) &= 2 \sum_{k=1}^{n_2} \lambda_{2,k}^2 + 4 \sum_{k=1}^{n_2} \lambda_{2,k}^2 \omega_k \\ &= 2 \text{tr}[(\mathbf{A}\boldsymbol{\Sigma})^2] + 4 \sum_{k=1}^{n_2} \lambda_{2,k}^2 \boldsymbol{\mu}'_y (\boldsymbol{\Phi}^+)' \mathbf{v}_{2,k} \mathbf{v}'_{2,k} \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y \\ &= 2 \text{tr}[(\mathbf{A}\boldsymbol{\Sigma})^2] + 4 \boldsymbol{\mu}'_y (\boldsymbol{\Phi}^+)' \sum_{k=1}^{n_2} \lambda_{2,k}^2 \mathbf{v}_{2,k} \mathbf{v}'_{2,k} \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y \\ &= 2 \text{tr}[(\mathbf{A}\boldsymbol{\Sigma})^2] + 4 \boldsymbol{\mu}'_y (\boldsymbol{\Phi}^+)' \mathbf{B}^2 \boldsymbol{\Phi}^+ \boldsymbol{\mu}_y \\ &= 2 \text{tr}[(\mathbf{A}\boldsymbol{\Sigma})^2] + 4 \boldsymbol{\mu}'_y (\boldsymbol{\Phi} \boldsymbol{\Phi}^+)' \mathbf{A} \boldsymbol{\Sigma} \mathbf{A} (\boldsymbol{\Phi} \boldsymbol{\Phi}^+) \boldsymbol{\mu}_y. \end{aligned} \quad (9.41)$$

If $\boldsymbol{\Sigma}$ has full rank of n , then $\mathcal{V}(q) = 2 \text{tr}[(\mathbf{A}\boldsymbol{\Sigma})^2] + 4 \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$, a very useful result.

Corollary 9.6.4 If Σ has full rank of n and $B = \Phi' A \Phi$ is idempotent ($\lambda_{2,k} = 1, \forall k$), then $n_2 = \text{rank}(B) = \text{rank}(A)$ and $q = \mathbf{y}' A \mathbf{y}$ is distributed as the sum of n_2 independent random variables distributed $\chi^2(1, \omega_k)$. Also $q \sim \chi^2(n_2, \mu_y' A \mu_y)$.

Proof. Partially left as an exercise. Lemma 1.30 is useful ($\Phi' A \Phi$ is idempotent $\Leftrightarrow \Sigma A$ is idempotent). The constituent matrix decomposition of a symmetric and idempotent matrix allows writing

$$\begin{aligned} \sum_{k=1}^{n_2} \omega_k &= \sum_{k=1}^{n_2} [\mu_y' (\Phi^+)' v_{2,k} v_{2,k}' \Phi^+ \mu_y] \\ &= \mu_y' \Phi^{-t} \left(\sum_{k=1}^{n_2} v_{2,k} v_{2,k}' \right) \Phi^{-1} \mu_y \\ &= \mu_y' \Phi^{-t} \left(V \begin{bmatrix} I_{n_2} & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{n-n_2} \end{bmatrix} V' \right) \Phi^{-1} \mu_y \\ &= \mu_y' \Phi^{-t} (\Phi' A \Phi) \Phi^{-1} \mu_y \\ &= \mu_y' A \mu_y. \end{aligned} \tag{9.42}$$

□

Theorem 9.7 If $\mathbf{y} \sim \mathcal{N}_n(\mu_y, \Sigma)$, $\text{rank}(\Sigma) = n$, and A_1 and A_2 are conforming constants, then

$$\mathcal{V}(\mathbf{y}' A_1 \mathbf{y}, \mathbf{y}' A_2 \mathbf{y}) = 2\text{tr}(A_1 \Sigma A_2 \Sigma) + 4\mu_y' A_1 \Sigma A_2 \mu_y. \tag{9.43}$$

Proof. Left as an exercise.

Theorem 9.8 If $q = \mathbf{y}' A \mathbf{y}$ with $\mathbf{y} \sim \mathcal{N}_n(\mu_y, \Sigma)$, $\text{rank}(\Sigma) = n$, $A = A'$ constant, and B ($m \times n$) constant (of any rank), then

$$\mathcal{V}(\mathbf{y}, \mathbf{y}' A \mathbf{y}) = 2\Sigma A \mu_y \tag{9.44}$$

$$\mathcal{V}(B \mathbf{y}, \mathbf{y}' A \mathbf{y}) = 2B \Sigma A \mu_y. \tag{9.45}$$

Proof. (Searle, 1971, p56) From the definition, the $n \times 1$ vector can be written

$$\begin{aligned} \mathcal{V}(B \mathbf{y}, \mathbf{y}' A \mathbf{y}) &= \text{E}\{ (B \mathbf{y} - B \mu_y) [\mathbf{y}' A \mathbf{y} - \mu_y' A \mu_y - \text{tr}(A \Sigma)] \} \\ &= \text{E}\{ B(\mathbf{y} - \mu_y) [(\mathbf{y} - \mu_y)' A (\mathbf{y} - \mu_y) + 2(\mathbf{y} - \mu_y)' A \mu_y - \text{tr}(A \Sigma)] \} \\ &= B \text{E}[(\mathbf{y} - \mu_y)(\mathbf{y} - \mu_y)' A (\mathbf{y} - \mu_y)] + 2B \text{E}[(\mathbf{y} - \mu_y)(\mathbf{y} - \mu_y)'] A \mu_y - \mathbf{0} \\ &= \mathbf{0} + (2B \Sigma A \mu_y) - \mathbf{0}, \end{aligned} \tag{9.46}$$

because the first and third moments of $\mathbf{y} - \mu_y$ are zero. □

In practice, the most frequent use of the following theorem is to infer from the fact $A = A^2$ that q must be chi square regardless of the true value of μ_y .

Theorem 9.9 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}_y, \mathbf{I}_n)$, $\mathbf{A} = \mathbf{A}'$ ($n \times n$), $\text{rank}(\mathbf{A}) = n_1$, and $\omega = \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$, then $q = \mathbf{y}' \mathbf{A} \mathbf{y} \sim \chi^2(n_1, \omega) \forall \boldsymbol{\mu}_y \Leftrightarrow \mathbf{A} = \mathbf{A}^2$.

Proof. (\Leftarrow) The spectral decomposition of \mathbf{A} is $\mathbf{A} = \mathbf{V} \text{Dg}(\lambda) \mathbf{V}'$ with $\mathbf{V}' \mathbf{V} = \mathbf{V} \mathbf{V}' = \mathbf{I}_n$. Given \mathbf{A} is idempotent we have

$$\mathbf{A} = [\mathbf{V}_1 \ \mathbf{V}_0] \begin{bmatrix} \mathbf{I}_{\nu} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{V}'_1 \\ \mathbf{V}'_0 \end{bmatrix} = \mathbf{V}_1 \mathbf{V}'_1. \tag{9.47}$$

Now $q = \mathbf{y}' \mathbf{A} \mathbf{y} = \mathbf{y}' \mathbf{V} \text{Dg}(\lambda) \mathbf{V}' \mathbf{y}$. If $\mathbf{z} = \mathbf{V}' \mathbf{y}$ then $\mathbf{y} = \mathbf{V} \mathbf{z}$,

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_0 \end{bmatrix} = \mathbf{V}' \mathbf{y} = \begin{bmatrix} \mathbf{V}'_1 \mathbf{y} \\ \mathbf{V}'_0 \mathbf{y} \end{bmatrix} \sim \mathcal{N}_n(\mathbf{V}' \boldsymbol{\mu}_y, \mathbf{I}_n), \tag{9.48}$$

and $q = \mathbf{y}' \mathbf{A} \mathbf{y} = \mathbf{z}' \text{Dg}(\lambda) \mathbf{z} = \mathbf{z}'_1 \mathbf{z}_1$. Having $\mathbf{z}_1 \sim \mathcal{N}_{n_1}(\mathbf{V}'_1 \boldsymbol{\mu}_y, \mathbf{I}_{\nu})$ implies $q = \mathbf{z}'_1 \mathbf{z}_1 \sim \chi^2(n_1, \omega)$ with $\omega = (\mathbf{V}'_1 \boldsymbol{\mu}_y)' (\mathbf{V}'_1 \boldsymbol{\mu}_y) = \boldsymbol{\mu}'_y \mathbf{V}_1 \mathbf{V}'_1 \boldsymbol{\mu}_y = \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$. \square

Proof. (\Rightarrow) Given $q = \mathbf{y}' \mathbf{A} \mathbf{y} \sim \chi^2(n_1, \omega)$ for every choice of $\boldsymbol{\mu}_y$, in which $\omega = \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$, and $n_1 = \text{rank}(\mathbf{A})$, we must prove \mathbf{A} is idempotent. The MGF of the distribution of $\mathbf{y}' \mathbf{A} \mathbf{y}$ is

$$m_{\mathbf{y}' \mathbf{A} \mathbf{y}}(t) = |\mathbf{I} - 2t\mathbf{A}|^{-1/2} \exp\left\{-\boldsymbol{\mu}'_y [\mathbf{I} - (\mathbf{I} - 2t\mathbf{A})^{-1}]^{-1} \boldsymbol{\mu}_y / 2\right\}, \tag{9.49}$$

while the MGF of $\chi^2(n_1, \omega)$ is

$$m_{\chi^2(n_1, \omega)}(t) = (1 - 2t)^{-\nu/2} \exp\left\{-(\omega/2)[1 - (1 - 2t)^{-1}]\right\}. \tag{9.50}$$

By assumption the two MGFs must be equal $\forall \boldsymbol{\mu}_y \in \mathfrak{R}^n$, including $\boldsymbol{\mu}_y = \mathbf{0}$. If $\boldsymbol{\mu}_y = \mathbf{0}$, then $\omega = 0$, which implies $(1 - 2t)^{-n_1/2} = |\mathbf{I} - 2t\mathbf{A}|^{-1/2} \forall t \in \mathfrak{R}$. For any (square) matrix $|\mathbf{I} - u\mathbf{A}| = \prod_{k=1}^n (1 - u\lambda_k)$, in which the eigenvalues of \mathbf{A} are $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$. If $u = 2t$, then raising both sides to the -2 power gives $(1 - u)^{n_1} = \prod_{k=1}^n (1 - u\lambda_k) \forall u \in \mathfrak{R}$. The two polynomials are identical and so must be of the same degree. The polynomial on the left is of degree $n_1 < n$, which implies the polynomial on the right is also of degree n_1 . In turn, the last $n - n_1$ terms are equal to 1 because $(1 - u\lambda_k) = 1$ for $k \in \{n_1+1, n_1+2, \dots, n\} \forall u \in \mathfrak{R}$ implies $\lambda_k = 0$ for $k \in \{n_1+1, n_1+2, \dots, n\}$. We now have $(1 - u)^{n_1} = \prod_{j=1}^{n_1} (1 - u\lambda_k) \forall u \in \mathfrak{R}$. Since the polynomials are identical, they must have identical roots. Therefore $\lambda_k = 1$ for $k \in \{1, 2, \dots, n_1\}$. The matrix $\mathbf{A} = \mathbf{A}'$ has n_1 eigenvalues of 1 and $n - n_1$ eigenvalues of 0. Finally, \mathbf{A} is idempotent of rank n_1 . \square

Corollary 9.9 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}_y, \boldsymbol{\Sigma})$, $\text{rank}(\boldsymbol{\Sigma}) = n$, $\mathbf{A} = \mathbf{A}'$ ($n \times n$), $\text{rank}(\mathbf{A}) = n_1$, and $\omega = \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$, then $\mathbf{y}' \mathbf{A} \mathbf{y} \sim \chi^2(n_1, \omega) \forall \boldsymbol{\mu}_y \Leftrightarrow (\mathbf{A}\boldsymbol{\Sigma}) = (\mathbf{A}\boldsymbol{\Sigma})^2$.

Proof. Factoring the covariance matrix allows transforming the random variables to independence. Spectral decomposition gives $\boldsymbol{\Sigma} = \mathbf{T} \text{Dg}(\lambda) \mathbf{T}' = \boldsymbol{\Phi} \boldsymbol{\Phi}'$

with $\Phi = \Upsilon \text{Dg}(\lambda)^{1/2}$. Alternately use the Cholesky decomposition of Σ , with Φ a lower triangular matrix. In general, Φ^{-1} exists if and only if Σ^{-1} exists. For $\text{rank}(\Sigma) = n$ we have $\Sigma^{-1} = \Phi^{-t} \Phi^{-1}$ and $\Phi^{-1} \Sigma \Phi^{-t} = I_n$. Also $\mathbf{z} = \Phi^{-1} \mathbf{y} \sim \mathcal{N}_n(\Phi^{-1} \boldsymbol{\mu}_y, I_n)$. If $\boldsymbol{\mu}_z = \Phi^{-1} \boldsymbol{\mu}_y$ then $\omega = \boldsymbol{\mu}'_z (\Phi' \mathbf{A} \Phi) \boldsymbol{\mu}_z = \boldsymbol{\mu}'_y \mathbf{A} \boldsymbol{\mu}_y$. Since $\mathbf{y} = \Phi \mathbf{z}$, we have $q = \mathbf{y}' \mathbf{A} \mathbf{y} = \mathbf{z}' \Phi' \mathbf{A} \Phi \mathbf{z}$. In terms of \mathbf{z} , $q \sim \chi^2(n_1, \omega)$ if and only if $\Phi' \mathbf{A} \Phi = (\Phi' \mathbf{A} \Phi)^2$. By Lemma 1.30 we conclude $q \sim \chi^2(n_1, \omega)$ if and only if $\mathbf{A} \Sigma$ is idempotent. \square

The immediately preceding results assumed a nonsingular Gaussian distribution, while Theorem 9.6 and the following theorem provide results for singular and nonsingular Gaussian distributions.

Theorem 9.10 For $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$ and $\text{rank}(\Sigma) = n_1 \leq n$, $q = (\mathbf{y}' \mathbf{A} \mathbf{y} + \mathbf{b}' \mathbf{y} + c)$. If $\mathbf{A} = \mathbf{A}'$ ($n \times n$), \mathbf{b} ($n \times 1$), c (1×1), $\nu = \text{tr}(\mathbf{A} \Sigma)$, and $\omega = (\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)' \Sigma (\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)$ are all fixed constants, then $q = (\mathbf{y}' \mathbf{A} \mathbf{y} + \mathbf{b}' \mathbf{y} + c) \sim \chi^2(\nu, \omega)$ if and only if all three of the following conditions hold regardless of the value of $\boldsymbol{\mu}$ and Σ :

1. $\Sigma \mathbf{A} \Sigma \mathbf{A} \Sigma = \Sigma \mathbf{A} \Sigma$
2. $(\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)' \Sigma = (\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)' \Sigma \mathbf{A} \Sigma$
3. $\boldsymbol{\mu}' \mathbf{A} \boldsymbol{\mu} + \mathbf{b}' \boldsymbol{\mu} + c = (\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)' \Sigma (\mathbf{A} \boldsymbol{\mu} + \mathbf{b}/2)$.

Searle (1971, Section 2.7) presented the theorem and its proof. One can obtain various corollaries by setting some of $\boldsymbol{\mu}$, \mathbf{b} , and c to zero.

The theorem is easily misinterpreted. If q has a noncentral chi-square distribution, the theorem does *not* imply $\mathbf{A} \Sigma$ is idempotent, only that all three of 1–3 are true. Searle (1971, p. 69) commented on problems which may arise from misinterpretation.

9.5 INDEPENDENCE AMONG LINEAR AND QUADRATIC FORMS

The proofs for Theorems 9.11 and 9.14 give independence properties based on factoring $\mathbf{A} = \mathbf{A}'$ as $\mathbf{A} = \mathbf{F} \mathbf{F}'$. If \mathbf{A} has any negative eigenvalues, then complex variables occur in \mathbf{F} . The complex variables could be avoided with a slight complication of the proofs. With $\mathbf{A} = \mathbf{V} \text{Dg}(\lambda) \mathbf{V}'$, if $s(\lambda_j) = 1$ for $\lambda_j \geq 0$ and -1 otherwise, writing $\text{Dg}(\lambda) = \text{Dg}(\{|\lambda_j|\}) \text{Dg}(\{s(\lambda_j)\})$ allows treating algebraic sign separately from eigenstructure.

Theorem 9.11 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$ with $\text{rank}(\Sigma) = n$, $n \times n$ constant $\mathbf{A} = \mathbf{A}'$ has rank $\nu \leq n$, and $m \times n$ \mathbf{B} is constant (and any rank), then $\mathbf{y}' \mathbf{A} \mathbf{y} \perp\!\!\!\perp \mathbf{B} \mathbf{y} \forall \boldsymbol{\mu} \Leftrightarrow \mathbf{B} \Sigma \mathbf{A} = \mathbf{0}$.

Proof. (\Leftarrow) Symmetric A can be written $A = FF'$ with F of dimension $n \times \nu$ and rank ν . One choice is $F = V_1 \text{Dg}(\lambda_1)^{1/2}$, with V_1 $n \times \nu$, and the columns the ν eigenvectors corresponding to the nonzero eigenvalues. Full column rank of F ensures $(F'F)^{-1}$ exists. Also

$$\begin{aligned} 0 &= B\Sigma A = B\Sigma FF' \\ &= B\Sigma FF'[F(F'F)^{-1}] \\ &= B\Sigma F \\ &= \mathcal{V}(By, F'y). \end{aligned} \tag{9.51}$$

The last line is true $\Leftrightarrow By \perp F'y'$ (due to the Gaussian distribution assumption), which implies $By \perp (F'y)'(F'y) \equiv y' Ay$.

Proof. (\Rightarrow) By Theorem 9.8 we know $\mathcal{V}(By, y' Ay) = 2B\Sigma A\mu$. By independence, $0 = 2B\Sigma A\mu \forall \mu$, which implies $B\Sigma A = 0$. \square

The first part of the proof reveals the underlying source of the independence, when it exists. Linear form By is independent of quadratic form $y' Ay$, because the quadratic form can be written as $y' Ay = y' FF'y = z'z$ with $z = F'y$ independent of By .

The theorem is both significant and has a familiar result as a special case. A standard result from univariate theory is that $y_i \sim \mathcal{N}(\mu, \sigma^2)$ i.i.d. $\Rightarrow \hat{\mu} = \bar{y}$ is independent of $\hat{\sigma}^2$. The following corollary states the familiar result formally.

Corollary 9.11 The sample mean is independent of the sample variance for i.i.d. Gaussian data. Given data $y \sim \mathcal{N}_N(\mu \mathbf{1}, \sigma^2 I_N)$, in the notation of the theorem, the usual estimators are $\bar{y} = b'y$ and $\hat{\sigma}^2 = y' Ay$ with $1 \times N$ $b' = (\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}'$ and $N \times N$ $A = [I - \mathbf{1}(\mathbf{1}'\mathbf{1})^{-1}\mathbf{1}']/(N - 1)$. Therefore \bar{y} and $\hat{\sigma}^2$ are independent.

Proof. It is easy to prove $b'(\sigma^2 I_N)A = 0$ because $\mathbf{1}'A = 0$, which allows applying the theorem. \square

The following theorem is the foundation for extending the preceding result to the singular Gaussian distribution.

Theorem 9.12 (Good, 1963) For $y \sim \mathcal{N}_n(0, \Sigma)$ with $\text{rank}(\Sigma) = n_1 \leq n$, constants $A = A'$ ($n \times n$), $a \in \mathbb{R}^n$, and $b \in \mathbb{R}^n$,
 (a) $y' Ay$ and $b'y$ are independent $\Leftrightarrow \Sigma A \Sigma b = 0$ and
 (b) $a'y$ and $b'y$ are independent $\Leftrightarrow a' \Sigma b = 0$.

Proof. Good (1963, Theorem 1C; corrigenda in 1966) provided a proof.

Theorem 9.13 If $y \sim \mathcal{N}_n(\mu, \Sigma)$, $\text{rank}(\Sigma) = n_1 \leq n$, $A = A'$ ($n \times n$) and B ($m \times n$) are constants, then $y' Ay$ and By are independent $\Leftrightarrow B\Sigma A \Sigma = 0$ and $B\Sigma A \mu = 0$.

Proof. (Searle, 1971, p. 70) Necessarily $\Sigma = \Phi\Phi'$, with Φ $n \times n_1$ and rank n_1 . Also $\mathbf{y} = \boldsymbol{\mu} + \Phi\mathbf{z}$, with $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I})$. If $\mathbf{b}' = \text{row}_i(\mathbf{B})$, then $\mathbf{y}'\mathbf{A}\mathbf{y} = \mathbf{z}'\Phi'\mathbf{A}\Phi\mathbf{z} + 2\boldsymbol{\mu}'\mathbf{A}\Phi\mathbf{z} + \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}$ and $\mathbf{b}'\mathbf{y} = \mathbf{b}'\Phi\mathbf{z} + \mathbf{b}'\boldsymbol{\mu}$. By Good's theorem, $\mathbf{z}'\Phi'\mathbf{A}\Phi\mathbf{z}$ is independent of $\mathbf{b}'\Phi\mathbf{z} \Leftrightarrow \mathbf{I}(\Phi'\mathbf{A}\Phi)\mathbf{I}(\Phi'\mathbf{b}) = \mathbf{0} \Leftrightarrow \mathbf{B}\Sigma\mathbf{A}\Sigma = \mathbf{0}$. Also $\boldsymbol{\mu}'\mathbf{A}\Phi\mathbf{z}$ is independent of $\mathbf{b}'\Phi\mathbf{z} \Leftrightarrow \boldsymbol{\mu}'\mathbf{A}\Phi\mathbf{I}\Phi'\mathbf{b} = \mathbf{0} \Leftrightarrow \mathbf{B}\Sigma\mathbf{A}\boldsymbol{\mu} = \mathbf{0}$. Combining results, $\mathbf{y}'\mathbf{A}\mathbf{y}$ is independent of $\mathbf{B}'\mathbf{y} \Leftrightarrow \mathbf{B}\Sigma\mathbf{A}\Sigma = \mathbf{0}$ and $\mathbf{B}\Sigma\mathbf{A}\boldsymbol{\mu} = \mathbf{0}$. \square

Corollary 9.13 Given the conditions of the theorem, if $\mathbf{B}\Sigma\mathbf{A} = \mathbf{0}$, then $\mathbf{y}'\mathbf{A}\mathbf{y}$ and $\mathbf{B}\mathbf{y}$ are independent.

Proof. Immediate from the theorem.

The useful corollary states a sufficient (but not necessary) condition for independence.

Theorem 9.14 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$ with $\text{rank}(\Sigma) = n$, $\mathbf{A} = \mathbf{A}'$ ($n \times n$) and $\mathbf{B} = \mathbf{B}'$ ($n \times n$) are constants, then $\mathbf{y}'\mathbf{A}\mathbf{y} \perp\!\!\!\perp \mathbf{y}'\mathbf{B}\mathbf{y} \forall \boldsymbol{\mu} \Leftrightarrow \mathbf{A}\Sigma\mathbf{B} = \mathbf{0}$.

Here $\mathbf{0} = \mathbf{A}\Sigma\mathbf{B} \Leftrightarrow \mathbf{0} = (\mathbf{A}\Sigma\mathbf{B})' = \mathbf{B}'\Sigma'\mathbf{A}' = \mathbf{B}\Sigma\mathbf{A}$. It is not necessary for either quadratic form to have a (marginal) chi-square distribution. The theorem is about independence only.

Proof. (\Leftarrow) The approach centers on defining a transformation and proving the independence of the two underlying linear forms. If $\mathbf{A} = \mathbf{F}_A\mathbf{F}_A'$ and $\mathbf{B} = \mathbf{F}_B\mathbf{F}_B'$, with \mathbf{F}_A and \mathbf{F}_B full (column) rank factors, then $(\mathbf{F}_A'\mathbf{F}_A)^{-1}$ and $(\mathbf{F}_B'\mathbf{F}_B)^{-1}$ exist. Furthermore

$$\begin{aligned} \mathbf{0} = \mathbf{A}\Sigma\mathbf{B} &= \mathbf{F}_A\mathbf{F}_A'\Sigma\mathbf{F}_B\mathbf{F}_B' \\ &= (\mathbf{F}_A'\mathbf{F}_A)^{-1}\mathbf{F}_A'\mathbf{F}_A\mathbf{F}_A'\Sigma\mathbf{F}_B\mathbf{F}_B'\mathbf{F}_B(\mathbf{F}_B'\mathbf{F}_B)^{-1} \\ &= \mathbf{F}_A'\Sigma\mathbf{F}_B \\ &= \mathcal{V}(\mathbf{F}_A'\mathbf{y}, \mathbf{F}_B'\mathbf{y}) = \mathcal{V}(\mathbf{x}, \mathbf{z}), \end{aligned} \quad (9.52)$$

with $\mathbf{x} = \mathbf{F}_A'\mathbf{y}$, $\mathbf{z} = \mathbf{F}_B'\mathbf{y}$. Under the assumption of Gaussian distribution, \mathbf{x} and \mathbf{z} are independent. Therefore $\mathbf{x}'\mathbf{x} = \mathbf{y}'\mathbf{A}\mathbf{y}$ is independent of $\mathbf{z}'\mathbf{z} = \mathbf{y}'\mathbf{B}\mathbf{y}$.

Proof. (\Rightarrow) Assuming independence of $\mathbf{y}'\mathbf{A}\mathbf{y}$ and $\mathbf{y}'\mathbf{B}\mathbf{y}$ allows writing

$$\begin{aligned} \mathcal{V}(\mathbf{y}'\mathbf{A}\mathbf{y} + \mathbf{y}'\mathbf{B}\mathbf{y}) &= \mathcal{V}(\mathbf{y}'\mathbf{A}\mathbf{y}) + \mathcal{V}(\mathbf{y}'\mathbf{B}\mathbf{y}) \\ &= \mathcal{V}[\mathbf{y}'(\mathbf{A} + \mathbf{B})\mathbf{y}]. \end{aligned} \quad (9.53)$$

The first part follows from the fact that $\mathcal{V}(\mathbf{y}'\mathbf{A}\mathbf{y}) + \mathcal{V}(\mathbf{y}'\mathbf{B}\mathbf{y}) - \mathcal{V}[\mathbf{y}'(\mathbf{A} + \mathbf{B})\mathbf{y}] = 0$. From a previous corollary,

$$\mathcal{V}(\mathbf{y}'\mathbf{A}\mathbf{y}) = 2\text{tr}[(\mathbf{A}\Sigma)^2] + 4\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{A}\boldsymbol{\mu} \quad (9.54)$$

$$\mathcal{V}(\mathbf{y}'\mathbf{B}\mathbf{y}) = 2\text{tr}[(\mathbf{B}\Sigma)^2] + 4\boldsymbol{\mu}'\mathbf{B}\Sigma\mathbf{B}\boldsymbol{\mu} \quad (9.55)$$

and

$$\begin{aligned} \mathcal{V}[\mathbf{y}'(\mathbf{A}+\mathbf{B})\mathbf{y}] &= 2\text{tr}\{[(\mathbf{A}+\mathbf{B})\Sigma]^2\} + 4\boldsymbol{\mu}'(\mathbf{A}+\mathbf{B})\Sigma(\mathbf{A}+\mathbf{B})\boldsymbol{\mu} \\ &= 2\text{tr}(\mathbf{A}\Sigma\mathbf{A}\Sigma) + 2\text{tr}(\mathbf{A}\Sigma\mathbf{B}\Sigma) + 2\text{tr}(\mathbf{B}\Sigma\mathbf{A}\Sigma) + \\ &\quad 2\text{tr}(\mathbf{B}\Sigma\mathbf{B}\Sigma) + 4\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{A}\boldsymbol{\mu} + 8\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu} + 4\boldsymbol{\mu}'\mathbf{B}\Sigma\mathbf{B}\boldsymbol{\mu}. \end{aligned} \quad (9.56)$$

The difference should be zero $\forall \boldsymbol{\mu}$ and Σ . Therefore

$$\begin{aligned} 0 &= 2\text{tr}(\mathbf{A}\Sigma\mathbf{B}\Sigma) + 2\text{tr}(\mathbf{B}\Sigma\mathbf{A}\Sigma) + 8\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu} \\ &= 4\text{tr}(\Sigma\mathbf{A}\Sigma\mathbf{B}) + 8\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu}. \end{aligned} \quad (9.57)$$

Letting $\boldsymbol{\mu} = \mathbf{0}$ implies $\text{tr}(\Sigma\mathbf{A}\Sigma\mathbf{B}) = 0$ which, with the equation above, implies $\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu} = 0 \forall \boldsymbol{\mu}$, which implies $\mathbf{A}\Sigma\mathbf{B} = \mathbf{0}$. \square

Theorem 9.15 If $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$ with $\text{rank}(\Sigma) = n_1 \leq n$, while $n \times n$ matrices \mathbf{A} and \mathbf{B} are constant (of any rank), then $\mathbf{y}'\mathbf{A}\mathbf{y}$ and $\mathbf{y}'\mathbf{B}\mathbf{y}$ are independent if and only if all three of the following hold:

- (1) $\boldsymbol{\mu}'\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu} = 0$, (2) $\Sigma\mathbf{A}\Sigma\mathbf{B}\boldsymbol{\mu} = \Sigma\mathbf{B}\Sigma\mathbf{A}\boldsymbol{\mu}$, and (3) $\Sigma\mathbf{A}\Sigma\mathbf{B}\Sigma = \mathbf{0}$.

Without loss of generality, \mathbf{A} and \mathbf{B} may be assumed symmetric.

Proof. Left as an exercise. The generalization for singular Gaussian vectors is from Searle (1971). The proof relies on Good's theorem.

Corollary 9.15.1 If $\mathbf{A}\Sigma\mathbf{B} = \mathbf{0}$, then $\mathbf{y}'\mathbf{A}\mathbf{y}$ and $\mathbf{y}'\mathbf{B}\mathbf{y}$ are independent.

The useful corollary states a sufficient (but not necessary) condition for independence.

Corollary 9.15.2 If both \mathbf{A} and \mathbf{B} are positive semidefinite or positive definite, then $\mathbf{y}'\mathbf{A}\mathbf{y}$ and $\mathbf{y}'\mathbf{B}\mathbf{y}$ are independent if and only if $\mathbf{A}\Sigma\mathbf{B} = \mathbf{0}$.

Proof. Shanbhag (1966; Searle, 1971, p. 71).

As one can see, relaxing the assumption of full-rank Σ to allow Σ to be positive semidefinite complicates the necessary conditions for independence. The basic tool for manipulating the singular Gaussian distribution is a transformation to a full-rank distribution, as follows. For $\mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \Sigma)$, $\text{rank}(\Sigma) = n_1 \leq n$, matrix Φ exists with full column rank such that $\Sigma = \Phi\Phi'$. In turn, $\mathbf{y} = \boldsymbol{\mu} + \Phi\mathbf{z}$ with $\mathbf{z} \sim \mathcal{N}_{n_1}(\mathbf{0}, \mathbf{I}_{n_1})$, and $\mathbf{z} = (\Phi'\Phi)^{-1}\Phi'(\mathbf{y} - \boldsymbol{\mu})$. The technique is central to the proof of the last theorem.

9.6 THE ANOVA THEOREM

ANOVA for general linear univariate models partitions the total sum of squares into component sums of squares as

$$SST = \mathbf{y}'\mathbf{y} = \mathbf{y}'\mathbf{A}_1\mathbf{y} + \mathbf{y}'\mathbf{A}_2\mathbf{y} + \cdots + \mathbf{y}'\mathbf{A}_n\mathbf{y}. \quad (9.58)$$

Typically, it is crucial to be able to assume the n sums of squares are totally independent and have marginal chi-square distributions. The ANOVA theorem provides necessary and sufficient conditions for such quadratic forms to have independent chi-square distributions. The theorem is founded on theory for idempotent matrices and has a geometric interpretation. If $\mathbf{A}_i = \mathbf{A}_i^2$ ($N \times N$), then $\mathbf{z} = \mathbf{A}_i\mathbf{y}$ ($N \times 1$) is a projection onto a subspace, and $\mathbf{z}'\mathbf{z} = \mathbf{y}'\mathbf{A}_i\mathbf{y}$ is the squared length of the vector. The theorem is about squared lengths of projections of \mathbf{y} ($N \times 1$) onto mutually orthogonal subspaces of the sample space.

The proof of the statistical theorem arises directly from a matrix theorem. In turn, Loynes' lemma simplifies proving the matrix theorem. We first prove Loynes' lemma, then derive the matrix decomposition used in the ANOVA theorem, and finally prove the ANOVA theorem itself. The presentation follows Searle (1971, p. 60–64). Early proofs were quite long. Banerjee (1964) produced a shorter proof, which was shortened and improved by Loynes (1966). Therefore Loynes' lemma is a key to a concise proof of the ANOVA theorem.

Lemma 9.1 (Loynes' lemma) If $\mathbf{B} = \mathbf{B}' = \mathbf{B}^2$, $\mathbf{Q} = \mathbf{Q}'$ is positive definite or positive semidefinite, and $\mathbf{I} - \mathbf{B} - \mathbf{Q}$ is positive definite or positive semidefinite, then $\mathbf{BQ} = \mathbf{0}$. Furthermore $\mathbf{BQ} = \mathbf{QB}$.

Proof. We prove $\mathbf{QBx} = \mathbf{0}$ for all $\mathbf{x} \in \mathfrak{R}^N$. It then follows $\mathbf{QB} = \mathbf{0}$. If $\mathbf{y} = \mathbf{Bx}$ for arbitrary $\mathbf{x} \in \mathfrak{R}^N$ then $\mathbf{y}'\mathbf{By} = \mathbf{y}'\mathbf{B}(\mathbf{Bx}) = \mathbf{y}'\mathbf{B}^2\mathbf{x} = \mathbf{y}'\mathbf{B}^1\mathbf{x} = \mathbf{y}'(\mathbf{Bx}) = \mathbf{y}'\mathbf{y}$. Thus $\mathbf{y}'(\mathbf{I} - \mathbf{B})\mathbf{y} = (\mathbf{y}'\mathbf{Iy} - \mathbf{y}'\mathbf{By}) = 0$ and $\mathbf{y}'(\mathbf{I} - \mathbf{B} - \mathbf{Q})\mathbf{y} = -\mathbf{y}'\mathbf{Qy} \geq 0$. The last inequality follows from the assumption $\mathbf{I} - \mathbf{B} - \mathbf{Q}$ is positive definite or positive semidefinite. The assumption \mathbf{Q} is positive definite or positive semidefinite implies by definition that $\mathbf{y}'\mathbf{Qy}$ is also n.n.d. for all \mathbf{y} . Hence $\mathbf{y}'\mathbf{Qy} = 0$. Since $\mathbf{Q} = \mathbf{Q}'$, \mathbf{F} exists such that $\mathbf{Q} = \mathbf{F}'\mathbf{F}$. Therefore $\mathbf{y}'\mathbf{F}'\mathbf{Fy} = 0$ implies $\mathbf{Fy} = \mathbf{0} = \mathbf{F}'\mathbf{Fy} = \mathbf{Qy} = \mathbf{Q}(\mathbf{Bx})$ for arbitrary \mathbf{x} . \square

Only \mathbf{B} is assumed to be idempotent. In many applications \mathbf{Q} is also idempotent. However, the lemma only requires the weaker condition of positive definite or positive semidefinite \mathbf{Q} . Of course, $\mathbf{Q} = \mathbf{Q}^2 \Rightarrow \mathbf{Q}$ positive definite or positive semidefinite. Similarly, $\mathbf{I} - \mathbf{B} - \mathbf{Q}$ is required to be only positive definite or positive semidefinite rather than idempotent. However, if either \mathbf{Q} or $(\mathbf{I} - \mathbf{B} - \mathbf{Q})$ is idempotent, then all of \mathbf{B} , \mathbf{Q} , $\mathbf{B} + \mathbf{Q}$, and $\mathbf{I} - \mathbf{B} - \mathbf{Q}$ are idempotent. Idempotency is easy to prove using the result $\mathbf{BQ} = \mathbf{0}$.

The matrices \mathbf{B} and \mathbf{Q} have the same eigenvectors. If \mathbf{A} and \mathbf{B} are symmetric matrices then $\mathbf{AB} = \mathbf{BA}$ if and only if \mathbf{A} and \mathbf{B} have the same eigenvectors

(Theorem 4.17, Schott, 2005). If so, $\mathbf{A} = \mathbf{V}_A \text{Dg}(\lambda_A) \mathbf{V}'_A$ and $\mathbf{B} = \mathbf{V}_B \text{Dg}(\lambda_B) \mathbf{V}'_B$. So $\mathbf{AB} = \mathbf{BA} = \mathbf{0}$ implies $\text{Dg}(\lambda_A) \text{Dg}(\lambda_B) = \mathbf{0}$.

Theorem 9.16 (Matrix decomposition used in the ANOVA theorem) If $\mathbf{A} = \mathbf{A}'$ ($N \times N$) of rank r is partitioned as $\mathbf{A} = \sum_{i=1}^k \mathbf{A}_i$ with each $\mathbf{A}_i = \mathbf{A}'_i$ of rank r_i , then the following conditions may be defined.

1. $\mathbf{A}_i = \mathbf{A}_i^2$ for $i \in \{1, 2, \dots, k\}$

2. $\mathbf{A}_i \mathbf{A}_{i'} = \mathbf{0} = \mathbf{A}_{i'} \mathbf{A}_i$ for $i \neq i'$

3. $\mathbf{A} = \mathbf{A}^2$

4. $r = \sum_{i=1}^k r_i$.

With the definitions, it follow that

I. Any two of 1, 2, 3 imply all of 1, 2, 3, and 4.

II. Together, 3 and 4 imply all of 1, 2, 3, and 4.

Proof. Proving I and II require only five steps:

1 and 3 \Rightarrow 2; 2 and 3 \Rightarrow 1; 1 and 2 \Rightarrow 3; 1 and 3 \Rightarrow 4; 3 and 4 \Rightarrow 1.

Proof that 1 and 3 \Rightarrow 2. By 3, $\mathbf{A} = \mathbf{A}^2 \Rightarrow (\mathbf{I} - \mathbf{A}) = (\mathbf{I} - \mathbf{A})^2 \Rightarrow (\mathbf{I} - \mathbf{A})$ is positive semidefinite (p.s.d.). By 1, $\mathbf{A}_k = \mathbf{A}_k^2 \Rightarrow \mathbf{A}_k$ is p.s.d. $\Rightarrow \sum_{i \neq k} \mathbf{A}_i$ is p.s.d. $\Rightarrow (\mathbf{A} - \mathbf{A}_k - \mathbf{A}_{i'})$ is p.s.d. Therefore $(\mathbf{I} - \mathbf{A}) + (\mathbf{A} - \mathbf{A}_k - \mathbf{A}_{i'}) = (\mathbf{I} - \mathbf{A}_k - \mathbf{A}_{i'})$ is p.s.d. By Loynes' lemma, $\mathbf{A}_k \mathbf{A}_{i'} = \mathbf{0}$, which implies 2.

Proof that 2 and 3 \Rightarrow 1. Eigenvector \mathbf{v} and corresponding eigenvalue λ of \mathbf{A}_i are defined by $\mathbf{A}_i \mathbf{v} = \mathbf{v} \lambda$ or $\lambda^{-1} \mathbf{A}_i \mathbf{v} = \mathbf{v}$ if $\lambda \neq 0$. By 2, $\mathbf{A}_k \mathbf{v} = \mathbf{A}_k \mathbf{A}_i \mathbf{v} / \lambda = \mathbf{0}$ for $k \neq i$ and $\lambda \neq 0$. For any nonzero eigenvalue of \mathbf{A}_i , $\lambda \neq 0$, and corresponding eigenvector \mathbf{v} , we have $\mathbf{A} \mathbf{v} = (\sum_k \mathbf{A}_k) \mathbf{v} = \mathbf{0} + \mathbf{A}_i \mathbf{v} = \lambda \mathbf{v}$. Therefore λ is an eigenvalue of \mathbf{A} . Since 3 $\Rightarrow \lambda = 1$ or 0, every nonzero eigenvalue of \mathbf{A}_i equals 1. Thus 1 holds.

Proof that 1 and 2 \Rightarrow 3. Using 1 and 2 in obvious ways we have 3, $\mathbf{A}^2 = \mathbf{A}$. Specifically $\mathbf{A}^2 = (\sum_k \mathbf{A}_k)^2 = \sum_i \sum_{i'} \mathbf{A}_i \mathbf{A}_{i'} = \sum_k \mathbf{A}_k^2 = \sum_k \mathbf{A}_k = \mathbf{A}$.

Proof that 1 and 3 \Rightarrow 4. Using 1 and 3 in obvious ways we have 4, $r = \sum_k r_k$, as follows: $r = \text{tr}(\mathbf{A}) = \text{tr}(\sum_k \mathbf{A}_k) = \sum_k \text{tr}(\mathbf{A}_k) = \sum_k r_k$.

Proof that 3 and 4 \Rightarrow 1. By 3, $\mathbf{A} = \mathbf{A}^2 \Rightarrow -(\mathbf{A} - \mathbf{I}) = (\mathbf{A} - \mathbf{I})^2 \Rightarrow \text{rank}(\mathbf{A} - \mathbf{I}) = N - r$. In turn, $(\mathbf{A} - \mathbf{I})$ has $N - r$ linearly independent columns and so $(\mathbf{A} - \mathbf{I}) \mathbf{x} = \mathbf{0}$ ($N \times 1$) is a set of N equations containing $N - r$ linearly independent (LIN) equations. Similarly, $-\mathbf{A}_i \mathbf{x} = \mathbf{0}$ is a set of N equations containing r_i LIN equations. The concatenation of n such sets of equations yields Nn equations,

$$\begin{bmatrix} \mathbf{A} - \mathbf{I} \\ -\mathbf{A}_2 \\ \vdots \\ -\mathbf{A}_n \end{bmatrix} \mathbf{x} = \mathbf{0}. \tag{9.59}$$

The equations contain at most $(N - r) + r_2 + r_3 + \dots + r_n$ LIN equations. By 4, a total of $N - r_1$ LIN equations exist. They can be reduced to N equations by adding terms on the left side to yield $(\mathbf{A}_1 - \mathbf{I})\mathbf{x} = \mathbf{0} \Leftrightarrow \mathbf{A}_1\mathbf{x} = \mathbf{x} \Leftrightarrow \mathbf{A}_1\mathbf{x} = \mathbf{1x}$. At most $N - r_1$ LIN equations exist.

Since the N equations $\mathbf{A}_1\mathbf{x} = \mathbf{1x}$ contain at most $N - r_1$ LIN equations, at least $N - (N - r_1) = r_1$ LIN solutions \mathbf{x} exist for the equations. Hence there are at least r_1 eigenvectors for \mathbf{A}_1 which correspond to eigenvalues equal to 1. Since $\text{rank}(\mathbf{A}_1) = r_1$, \mathbf{A}_1 only has r_1 nonzero eigenvalues. Therefore $\lambda = 1$ is an eigenvalue of multiplicity r_1 and $\lambda = 0$ is an eigenvalue of multiplicity $N - r_1$. Consequently \mathbf{A}_1 must be idempotent. The same proof can be applied to any other \mathbf{A}_i . Thus, by extension, all \mathbf{A}_i must be idempotent, which gives 1. \square

Theorem 9.17 (The ANOVA theorem) If $\mathbf{y} \sim \mathcal{N}_N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ with $\text{rank}(\boldsymbol{\Sigma}) = N$, $\mathbf{A} = \mathbf{A}'$ ($N \times N$) of rank r with $\mathbf{A} = \sum_{i=1}^k \mathbf{A}_i$, $\mathbf{A}_i = \mathbf{A}'_i$ with rank r_i , and $q = \mathbf{y}'\mathbf{A}\mathbf{y}$, the following conditions may be defined.

1. $\mathbf{A}_i\boldsymbol{\Sigma} = (\mathbf{A}_i\boldsymbol{\Sigma})^2$ for $i \in \{1, 2, \dots, k\}$, which is equivalent to $\mathbf{A}_i\boldsymbol{\Sigma}\mathbf{A}_i = \mathbf{A}_i$,
2. $\mathbf{A}_i\boldsymbol{\Sigma}\mathbf{A}_{i'} = \mathbf{0}$ for all $i \neq i'$,
3. $\mathbf{A}\boldsymbol{\Sigma} = (\mathbf{A}\boldsymbol{\Sigma})^2$, and
4. $r = \sum_{i=1}^k r_i$.

For the conditions defined,

- (a) $\mathbf{y}'\mathbf{A}_i\mathbf{y} \sim \chi^2(r_i, \boldsymbol{\mu}'\mathbf{A}_i\boldsymbol{\mu})$,
- (b) $\mathbf{y}'\mathbf{A}_i\mathbf{y}$ is independent of $\mathbf{y}'\mathbf{A}_{i'}\mathbf{y} \ \forall i \neq i'$, and
- (c) $\mathbf{y}'\mathbf{A}\mathbf{y} \sim \chi^2(r, \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu})$

are all simultaneously true if and only if

- I. any two of 1, 2, and 3 are true, or
- II. 3 and 4 are both true.

Corollary 9.17.1 (Cochran, 1934) If $\mathbf{y} \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}_N)$ with $\mathbf{A} = \mathbf{I}_N$ [and $\text{rank}(\mathbf{A}) = N$] is partitioned as $\mathbf{I} = \sum_{i=1}^k \mathbf{A}_i$ with $\mathbf{A}_i = \mathbf{A}'_i$ of rank r_i , and $q = \mathbf{y}'\mathbf{A}\mathbf{y}$, then the variates $q_i = \mathbf{y}'\mathbf{A}_i\mathbf{y}$, $i \in \{1, 2, \dots, k\}$, are mutually independent and distributed as $\chi^2(r_i, 0) \Leftrightarrow N = \sum_{i=1}^k r_i$. Obviously $q \sim \chi^2(N, 0)$.

Proof. Left as an exercise.

Example 9.4 A $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$ with Gaussian distribution has

$$\{W_1, W_2, W_3, W_4\} = \left\{ \begin{matrix} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, & \begin{bmatrix} -1 \\ -1 \\ -1 \\ 1 \\ 1 \end{bmatrix}, & \begin{bmatrix} -1 & 1 \\ 0 & -2 \\ 1 & 1 \\ -1 & 1 \\ 0 & -2 \\ 1 & 1 \end{bmatrix}, & \begin{bmatrix} 1 & -1 \\ 0 & 2 \\ -1 & -1 \\ -1 & 1 \\ 0 & -2 \\ 1 & 1 \end{bmatrix} \end{matrix} \right\}. \quad (9.60)$$

If $X = [W_1 \ W_2 \ W_3]$ and $\beta' = [\beta'_1 \ \beta'_2 \ \beta'_3]$, then $(X'X) = Dg(\{6, 6, 4, 12\})$ and $(X'X)^{-1} = Dg(\{1/6, 1/6, 1/4, 1/12\})$. Finally, if

$$A_k = \sigma^{-2}W_k(W'_k W_k)^{-1}W'_k, \quad (9.61)$$

then

$$A = \sigma^{-2}I = \sum_{k=1}^4 A_k. \quad (9.62)$$

The $\{A_k\}$ provide projections of y ($N \times 1$) onto mutually orthogonal subspaces of the sample space. The subspace spanned by $(A_1 + A_2 + A_3)$ (with basis X) is the *estimation space*. The subspace spanned by A_4 (with basis W_4) is the *error space*. The corresponding source table is given below. The presence of σ^{-2} in the definition of A_k scales the underlying Gaussian variables to have unit variance, which leads the associated sums of squares to be chi square, as required in the preceding corollary.

ANOVA Table for Example

Source	df	Sum of Squares
Mean	1	$y'A_1y$
A	1	$y'A_2y$
B	2	$y'A_3y$
Residual	2	$y'A_4y$
Total	6	$y'Ay = y'y\sigma^{-2}$

Corollary 9.17.2 In a $GLM_{N,q}FR(y_i; X_i\beta, \sigma^2)$ with Gaussian errors, the fixed and known constants C and θ_0 define the a priori secondary parameter ($a \times 1$) $\theta = C\beta - \theta_0$, with $\text{rank}(C) = a \leq q$. In turn,

$$F(y) = \frac{(\hat{\theta} - \theta_0)'[C(X'X)^{-1}C']^{-1}(\hat{\theta} - \theta_0)/a}{\hat{\sigma}^2} = \frac{SSH/a}{SSE/(N - q)} \quad (9.63)$$

is a ratio of quadratic forms. Here $SSH =$ sum of squares for the hypothesis is computed from $\hat{\beta} = (X'X)^{-1}X'y$ and $\hat{\theta} = C\hat{\beta}$, while

$$\begin{aligned}
 SSE &= \text{sum of squares for error} \\
 &= (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \\
 &= \mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}.
 \end{aligned}
 \tag{9.64}$$

Hence $F(\mathbf{y}) \sim F(a, N - q, \omega)$ for $\omega = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)/\sigma^2$.

Proof. Since $\hat{\boldsymbol{\beta}}$ and $(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})$ are statistically independent, SSE and SSH are statistically independent. Furthermore we can prove $SSE/\sigma^2 \sim \chi^2(N - q)$ and $SSH/\sigma^2 \sim \chi^2(a, \omega)$. Thus $F(\mathbf{y})$ is distributed as a ratio of independent chi-square statistics divided by their degrees of freedom. \square

We shall see in Chapter 15 that in testing $H_0 = \mathbf{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbf{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$ both the likelihood ratio test (LRT) procedure and the union intersection test (UIT) procedure can use $F(\mathbf{y})$ as the test statistic.

9.7 RATIOS INVOLVING QUADRATIC FORMS

The F and t distributions play central roles in describing distributions of test statistics for linear models. Both are defined in terms of other random variables. In particular, if $z \sim \mathcal{N}(0, 1)$ and $x \sim \chi^2(\nu)$ are independent, then

$$t = \frac{z + \mu}{\sqrt{x/\nu}}.
 \tag{9.65}$$

The situation will be indicated by writing $t \sim t(\nu, \mu)$ for the noncentral case or $t \sim t(\nu)$ for the central case. The corresponding CDF is indicated $F_t(t_*, \nu, \mu)$, and the corresponding quantile is $F_t^{-1}(p; \nu, \mu)$. A two-tailed test of size α use $t_{\text{crit}} = F_t^{-1}(1 - \alpha/2; \nu)$, while a one tailed test uses either $F_t^{-1}(1 - \alpha; \nu)$ or $F_t^{-1}(\alpha; \nu)$, depending on the direction (sign of t) required.

If $x_1 \sim \chi^2(\nu_1, \omega_1)$ is independent of $x_2 \sim \chi^2(\nu_2, \omega_2)$, then

$$f = \frac{x_1/\nu_1}{x_2/\nu_2}
 \tag{9.66}$$

is described as following a doubly noncentral F distribution, $f \sim F(\nu_1, \nu_2, \omega_1, \omega_2)$. The (singly) noncentral F has $\omega_2 = 0$, written $f \sim F(\nu_1, \nu_2, \omega_1)$, and the central has $\omega_1 = \omega_2 = 0$, written $f \sim F(\nu_1, \nu_2)$. Noncentral F has corresponding CDF indicated $F_F(f; \nu_1, \nu_2, \omega_1)$ and quantile $f_{\text{crit}} = F_F^{-1}(p; \nu_1, \nu_2, \omega_1)$. Equivalently $t^2 \sim F(1, \nu_2, \omega_1)$. A size α test uses $f_{\text{crit}} = F_F^{-1}(1 - \alpha; \nu_1, \nu_2)$, which corresponds to a two-tailed t test if $\nu_1 = 1$. A one-tailed t test may be performed with $F_F^{-1}(1 - 2\alpha; 1, \nu_2)$ while requiring the underlying t to have the correct sign. The notation may be summarized by writing $F_F(f_0; \nu_1, \nu_2, \omega) = \Pr\{f \leq f_0\}$ and, for $\omega = 0$, $\Pr\{f > f_{\text{crit}}\} = \alpha$.

More general ratios of quadratic forms occur naturally in linear models. Such a variable can be written

$$\frac{q_1}{q_2} = \frac{\sum_{j=1}^{J_1} c_j x_j}{\sum_{j=J_1+1}^{J_1+J_2} c_j x_j} \tag{9.67}$$

for constants $\{c_j\}$ and random $x_j \sim \chi^2(\nu_j, \omega_j)$ all fully independent. Typically $c_j > 0$, which implies $\Pr\{q_1/q_2 > 0\} = 1$. For $r_0 > 0$ a simple transformation gives

$$\begin{aligned} \Pr\{q_1/q_2 \leq r_0\} &= \Pr\left\{\sum_{j=1}^{J_1} c_j x_j - r_0 \sum_{j=J_1+1}^{J_1+J_2} c_j x_j \leq 0\right\} \\ &= \Pr\left\{\sum_{j=1}^{J_1+J_2} d_j x_j \leq 0\right\}. \end{aligned} \tag{9.68}$$

Although $s = \sum_{j=1}^{J_1+J_2} d_j x_j$ has a simple and known characteristic function, computing $\Pr\{s \leq 0\} = \Pr\{q_1/q_2 \leq r_0\}$ proves difficult. With $\{d_j\}$ known, Davies' (1980) algorithm allows computing $\Pr\{s \leq 0\}$ and more general results with specifiable precision. The method uses numerical inversion of the characteristic function (a numerical integration). Interest in special cases had led to the development of many alternative algorithms, typically based on series expansions. Johnson and Kotz (1970, Chapter 29, p. 169–173) summarized many issues of theory and computational practice in the long history of the problem. Johnson, Kotz, and Balakrishnan (1994) provided some additional information.

In many settings, an approximation provides sufficient accuracy and can be much faster to compute. For $c_j > 0$ and $\omega_j \equiv 0$, a Satterthwaite (1946) approximation matches the mean and variance of q_k to q_{*k} , with $q_{*k}/\lambda_{*k} \sim \chi^2(\nu_{*k}, 0)$. In turn, $\Pr\{q_1/q_2 \leq r_0\} \approx F_F[r_0(\lambda_{*2}\nu_{*2})/(\lambda_{*k}\nu_{*1}); \nu_{*1}, \nu_{*2}]$. Kim, Gribbin, Muller and Taylor (2005) generalized the approximation to allow $\omega_j > 0$ for $j \leq J_1$ by using $q_{*1}/\lambda_{*1} \sim \chi^2(\nu_{*1}, \omega_{*1})$.

EXERCISES

9.1 Suppose $\mathbf{y} \sim \mathcal{N}_n(\mathbf{1}_n\mu, \Sigma)$ and $\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho & \cdot & \rho \\ \rho & 1 & \cdot & \rho \\ \cdot & \cdot & \cdot & \cdot \\ \rho & \rho & \cdot & 1 \end{bmatrix}$. Thus $E(y_i) = \mu$ for all

i , $\mathcal{V}(y_i) = \sigma^2$ for all i , and $\mathcal{V}(y_i, y_j) = \sigma^2\rho$ for all $i \neq j$; that is, the y 's are equicorrelated. Equivalently, $\Sigma = \sigma^2[(1 - \rho)\mathbf{I} + \rho\mathbf{1}\mathbf{1}']$.

9.1.1 Show that $\sum_{i=1}^n (y_i - \bar{y})^2 / [\sigma^2(1 - \rho)]$ is $\chi^2(n - 1)$.

9.1.2 Given that \mathbf{V}_T is an $n \times (n - 1)$ matrix which is columnwise orthonormal such that $\mathbf{V}_T'\mathbf{1}_n = \mathbf{0}$, find the distribution of $\mathbf{y}_T = \mathbf{V}_T'\mathbf{y}$.

9.1.3 Explicitly specify the distribution of $Q = \mathbf{y}'_T\mathbf{y}_T$.

9.2 Suppose $\mathbf{y} \sim \mathcal{N}_3(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, with $\boldsymbol{\mu} = \begin{bmatrix} 2 \\ -1 \\ 3 \end{bmatrix}$, $\boldsymbol{\Sigma} = \begin{bmatrix} 4 & 1 & 0 \\ 1 & 2 & 1 \\ 0 & 1 & 3 \end{bmatrix}$, and

$$\mathbf{A} = \begin{bmatrix} 1 & -3 & -8 \\ -3 & 2 & -6 \\ -8 & -6 & 3 \end{bmatrix}.$$

9.2.1 Find $E(\mathbf{y}'\mathbf{A}\mathbf{y})$.

9.2.2 Find $\mathcal{V}(\mathbf{y}'\mathbf{A}\mathbf{y})$.

9.2.3 Fully specify the exact distribution of $\mathbf{y}'\mathbf{A}\mathbf{y}$ (including all parameters and dimensions). Explain why the distribution you claim provides the correct answer.

You may use IML or any other matrix language for numerical calculations.

If you do, please show both the code and final results to help the grader.

9.2.4 If $\boldsymbol{\Sigma} = \sigma^2\mathbf{I}$, does $\mathbf{y}'\mathbf{A}\mathbf{y}/\sigma^2$ have a χ^2 distribution?

9.3 Assume $x_1/\sigma^2 \sim \chi^2(\nu_1, \omega)$, independent of x_2 , and $x_2/\sigma^2 \sim \chi^2(\nu_2)$.

Define $r = (x_1/\nu_1)/(x_2/\nu_2)$, $c = \nu_2/(\nu_1 r_0)$ for $r_0 > 0$, and $s = cx_1 - x_2$.

9.3.1 Prove that $\Pr\{r \leq r_0\} = \Pr\{s \leq 0\}$.

9.3.2 What is $E(s)$?

9.3.3 What is $\mathcal{V}(s)$?

9.3.4 What is the moment generating function of s ?

CHAPTER 10

Multivariate Quadratic Forms

10.1 THE WISHART DISTRIBUTION

With \mathbf{A} constant and \mathbf{y} Gaussian, $q = \mathbf{y}'\mathbf{A}\mathbf{y}$ is a univariate quadratic form and equals a weighted sum of chi squares. With equal weights, q is a scaled chi square. The necessary and sufficient conditions are detailed in the previous chapter. Replacing vector Gaussian \mathbf{y} by matrix Gaussian leads to the following definition.

Definition 10.1 The $N \times N$ constant $\mathbf{A} = \mathbf{A}'$ and $\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{\Xi}, \Sigma)$ create the $p \times p$ matrix $\mathbf{Q} = \mathbf{Y}'\mathbf{A}\mathbf{Y}$, a *multivariate quadratic form*.

In the multivariate case we must always assume $\mathbf{A} = \mathbf{A}'$. Idempotent \mathbf{A} leads to a special distribution (chi square) in the univariate case. Similarly, idempotent \mathbf{A} and $\mathbf{\Xi}$ leads to \mathbf{Q} having a Wishart distribution, which is one multivariate generalization (of many) of the chi square. Wishart (1928), Johnson and Kotz (1972), Arnold (1981), Muirhead (1984), Gupta and Nagar (2000), and Anderson (2004) included related treatments. The nomenclature and approach to many proofs used here closely follow the presentation in Muller and Chi (2006).

Definition 10.2 (a) If $\mathbf{Y} \sim \mathcal{N}_{\nu,p}(\mathbf{0}, \mathbf{I}_\nu, \Sigma)$, then $\mathbf{Y}'\mathbf{Y} \sim \mathcal{W}_p(\nu, \Sigma)$ indicates $\mathbf{Y}'\mathbf{Y}$ follows a *central (integer) Wishart* distribution with (integer) $\nu > 0$ degrees of freedom.

(b) If $\mathbf{Y} \sim \mathcal{N}_{\nu,p}(\mathbf{M}, \mathbf{I}_\nu, \Sigma)$, then $\mathbf{Y}'\mathbf{Y} \sim \mathcal{W}_p(\nu, \Sigma, \mathbf{M}'\mathbf{M})$ indicates $\mathbf{Y}'\mathbf{Y}$ follows a *noncentral (integer) Wishart* distribution with (integer) $\nu > 0$ degrees of freedom, *shift* $\mathbf{\Delta} = \mathbf{M}'\mathbf{M}$, and *noncentrality* $\mathbf{\Omega} = \mathbf{M}'\mathbf{M}\Sigma^+$.

(c) Singular Σ may be emphasized by writing $\mathcal{SW}_p(\nu, \Sigma)$ or $\mathcal{SW}_p(\nu, \Sigma, \mathbf{\Delta})$.

(d) Writing $(\mathcal{S})\mathcal{W}_p(\nu, \Sigma)$ or $(\mathcal{S})\mathcal{W}_p(\nu, \Sigma, \mathbf{\Delta})$ indicates possibly singular Σ .

In parallel to a chi square, a Wishart is defined as a quadratic form of independent standard Gaussian random vectors. The Wishart definition also includes scale. The shift and noncentrality parameters reflect the nature of the underlying Gaussian variables. If (and only if) $\mathbf{M} = \mathbf{0}$, then $\mathbf{\Delta} = \mathbf{M}'\mathbf{M} = \mathbf{0}$ and $\mathbf{\Omega} = \mathbf{0}$, which reduces a noncentral Wishart to a central. Eigenvalues of the

noncentrality matrix Ω are invariant to any full-rank transformation of the columns of \mathbf{Y} . Therefore they are scale free, in the sense of being invariant to multiplying each variable by a possibly distinct nonzero constant.

Historically, most discussions of the Wishart have assumed full-rank Σ , which implies $\Omega = \mathbf{M}'\mathbf{M}\Sigma^{-1}$ is unique, and led to the notation $\mathcal{W}_p(\nu, \Sigma, \Omega)$. For singular Σ , some authors define noncentrality as $\mathbf{M}'\mathbf{M}\Sigma^-$ or $\mathbf{M}'\mathbf{M}\Sigma^+$. In many applications, only functions of $\mathbf{M}'\mathbf{M}\Sigma^-$ invariant to the choice of Σ^- occur, which allows using $\mathbf{M}'\mathbf{M}\Sigma^+$ without loss of generality.

Defining the Wishart in terms of $\Delta = \mathbf{M}'\mathbf{M}$ not only avoids the ambiguity in singular cases but also is consistent with chi-square notation. If $z \sim \mathcal{N}(\mu_z, 1)$, then $y = z\sigma \sim \mathcal{N}(\mu_y, \sigma^2)$ with $\mu_y = \sigma\mu_z$. In turn, $y^2/\sigma^2 = z^2 \sim \chi^2(1, \mu_z^2)$ and $z^2 \sim \mathcal{W}_1(1, \mu_z^2)$, with noncentrality $\mu_z^2 = \mu_y^2/\sigma^2 = \omega$, while $y^2 \sim \mathcal{W}_1(1, \mu_y^2)$. If $\mathbf{Z} \sim \mathcal{N}_{\nu,p}(\mathbf{M}_Z, \mathbf{I}_\nu, \mathbf{I}_p)$ and $\Sigma = \Phi\Phi'$, then $\mathbf{Y} = \mathbf{Z}\Phi' \sim \mathcal{N}_{\nu,p}(\mathbf{M}_Y, \mathbf{I}_\nu, \Sigma)$ with $\mathbf{M}_Y = \mathbf{M}_Z\Phi'$. In turn, $\Phi^{-1}\mathbf{Y}'\mathbf{Y}\Phi^{-t} = \mathbf{Z}'\mathbf{Z} \sim \mathcal{W}_p(\nu, \mathbf{I}_p, \mathbf{M}'_Z\mathbf{M}_Z)$ with noncentrality $\mathbf{M}'_Z\mathbf{M}_Z = \Phi^{-1}\mathbf{M}'_Y\mathbf{M}_Y\Phi^{-t} = \Omega$, while $\mathbf{Y}'\mathbf{Y} \sim \mathcal{W}_p(\nu, \Sigma, \mathbf{M}'_Y\mathbf{M}_Y)$.

Table 10.1 Impact of Rank Conditions on Eigenvalue Estimation:

Central Wishart, $\mathbf{S} = \nu\widehat{\Sigma} \sim \mathcal{W}_p(\nu, \Sigma)$

	Σ Singular $\text{rank}(\Sigma) = p_1 < p$	Σ Nonsingular $\text{rank}(\Sigma) = p_1 = p$	$\forall \lambda_j \neq 0$ Estimable?
$\text{rank}(\widehat{\Sigma}) = \nu < p_1$	$0 < \nu < p_1$ $< p$	$0 < \nu < p_1 = p$	No
$\text{rank}(\widehat{\Sigma}) = p_1 \leq \nu$	0 $< p_1$ $< p \leq \nu$ 0 $< p_1 \leq \nu < p$	0 $< p_1 = p \leq \nu$	Yes Yes

Arnold (1981, p. 317) noted that various authors have used the term “singular” or “pseudo” if Σ is singular or if $\nu < p$. Such approaches fail to describe all possible combinations in Table 10.1. We suggest the following terms. The distinction between *population singular* and *population nonsingular* specifies $\text{rank}(\Sigma) = p$ or $\text{rank}(\Sigma) = p_1 < p$. In turn, $\widehat{\Sigma}$ may be singular due to $\nu < p$ or $\text{rank}(\Sigma) = p_1 < p$. The contrast between *sample-rank sufficient* (to estimate all nonzero population eigenvalues) and *sample-rank insufficient* fully captures the necessary distinction. Estimated eigenvalues are roots of the scalar polynomial $|\widehat{\Sigma} - \widehat{\lambda}\mathbf{I}_p| = 0$. Hence $\text{rank}(\widehat{\Sigma})$ determines the number of eigenvalues that can be estimated. As long as $\text{rank}(\widehat{\Sigma}) = \text{rank}(\Sigma)$ all population eigenvalues can be estimated, which reflects whether $\nu \geq p_1$ or $\nu < p_1$. Although four of five cases have singular $\widehat{\Sigma}$ in Table 10.1, only two of five are sample-rank insufficient.

Since $\mathbf{S} = \mathbf{S}'$, only $p(p + 1)/2$ distinct elements exist, and the p^2 elements are not functionally independent. Saying “ \mathbf{S} follows a Wishart distribution” always refers to the distribution of $\mathbf{z} = \text{vech}(\mathbf{S})$. Only Wishart matrices that are both population and sample nonsingular have a density. In parallel to the singular Gaussian, a population-nonsingular Wishart can be extracted from a population-

singular one to represent all of the information available. The statement holds whether or not the Wishart is sample-rank sufficient or insufficient. In contrast, sample-rank insufficiency can *not* be converted to sample-rank-sufficiency, except for special cases of simple covariance structure.

Example 10.1 A professor interested in predicting Graduate Record Examination (GRE) scores from undergraduate grade point might examine the residual covariance matrix for the scores on Verbal, Quantitative and Analytic sections of the test. The test construction process makes it reasonable to assume the data are multivariate Gaussian. With data from $N = 20$ students, the covariance estimate $\widehat{\Sigma}_1$ would be such that $\mathbf{S}_1 = 19\widehat{\Sigma}_1 \sim \mathcal{W}_3(20 - 1, \Sigma_1)$, which is population nonsingular and sample-rank sufficient.

If the professor includes Total = Verbal + Quantitative + Analytic, then $\mathbf{S}_2 = 19\widehat{\Sigma}_2 \sim \mathcal{SW}_4(20 - 1, \Sigma_2)$, which is population-singular and sample-rank-sufficient. Infinitely many 4×3 transformation matrices can transform \mathbf{S}_2 into a population-nonsingular Wishart. Choices include

$$\mathbf{T}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} = \begin{bmatrix} \mathbf{I}_3 \\ \mathbf{0} \end{bmatrix} \tag{10.69}$$

$$\mathbf{T}_2 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \tag{10.70}$$

As will be proven later in the chapter, $\mathbf{T}'_j \mathbf{S}_2 \mathbf{T}_j \sim \mathcal{W}_3(20 - 1, \mathbf{T}'_j \Sigma_2 \mathbf{T}_j)$, which corresponds to studying $\mathbf{Y}\mathbf{T}_j$ with dimensions $(20 \times 4)(4 \times 3)$.

Example 10.2 Medical imaging and genetic scientists often encounter data with more variables than participants. Pizer et al. (2003) compared human and computer segmentations from CT images of $N = 12$ kidneys. They modeled the surfaces at 88 points, giving $p = 88 \cdot 3 = 264$ (x, y, z) variables, describing location in three dimensions. Gaussian data lead to the assumption of an observed covariance matrix $\widehat{\Sigma}$ such that $\mathbf{S} = 11\widehat{\Sigma} \sim \mathcal{W}_{264}(12 - 1, \Sigma)$. Here \mathbf{S} is population nonsingular and sample-rank insufficient.

10.2 THE CHARACTERISTIC FUNCTION OF THE WISHART

The following lemma summarizes properties of any covariance matrix. The notation helps describe the characteristic function of the Wishart.

Lemma 10.1 A $p \times p$ covariance matrix is symmetric with no negative eigenvalues, $\{\lambda_j\}$, which allows writing $\Sigma = \Upsilon \text{Dg}(\lambda) \Upsilon' = \Phi \Phi'$, with Υ , $\text{Dg}(\lambda)$, and $\Phi = \Upsilon \text{Dg}(\lambda)^{1/2}$ all $p \times p$. Considering p_1 columns in the first of

two partitions, with orthonormal Υ , gives $\Upsilon = [\Upsilon_1 \ \Upsilon_0]$, $\lambda = [\lambda_1' \ \lambda_0']'$, $\Phi = [\Phi_1 \ \Phi_0]$, $\Phi_1 = \Upsilon_1 \text{Dg}(\lambda_1)^{1/2}$, $\Phi_0 = \Upsilon_0 \text{Dg}(\lambda_0)^{1/2}$, and

$$\begin{aligned} \Sigma &= [\Upsilon_1 \ \Upsilon_0] \begin{bmatrix} \text{Dg}(\lambda_1) & \mathbf{0} \\ \mathbf{0} & \text{Dg}(\lambda_0) \end{bmatrix} \begin{bmatrix} \Upsilon_1' \\ \Upsilon_0' \end{bmatrix} \\ &= [\Phi_1 \ \Phi_0] \begin{bmatrix} \Phi_1' \\ \Phi_0' \end{bmatrix} \\ &= \Phi_1 \Phi_1' + \Phi_0 \Phi_0'. \end{aligned} \tag{10.1}$$

If $\text{rank}(\Sigma) = p_1 \leq p$, then $\Sigma = \Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1' = \Phi_1 \Phi_1'$ and $\Phi_1^+ = (\Phi_1' \Phi_1)^{-1} \Phi_1' = \text{Dg}(\lambda_1)^{-1/2} \Upsilon_1'$.

Proof. Eigenanalysis and partitioned matrix properties give the results.

Theorem 10.1 (a) For $\nu > 0$, $\Sigma = \Phi \Phi' = \Phi_1 \Phi_1'$, $p \times p \ \Phi$, $p \times p_1 \ \Phi_1$, $\text{rank}(\Phi) = \text{rank}(\Phi_1) = \text{rank}(\Sigma) = p_1 \leq p$, $i = (-1)^{1/2}$, $p \times p$ real $U = U'$, $\langle T \rangle_{jj} = u_{jj}$, and $\langle T \rangle_{jk} = u_{jk}/2$, the characteristic function of $S \sim (S) \mathcal{W}_p(\nu, \Sigma)$ is

$$\phi_S(T) = |I_p - 2iT\Sigma|^{-\nu/2} \tag{10.2}$$

$$= |I_p - 2iT\Phi\Phi'|^{-\nu/2} \tag{10.3}$$

$$= |I_p - 2iT\Phi_1\Phi_1'|^{-\nu/2} \tag{10.4}$$

$$= |I_p - 2i\Phi_1' T \Phi_1|^{-\nu/2} \tag{10.5}$$

$$= |I_{p_1} - 2i\Phi_1' T \Phi_1|^{-\nu/2}. \tag{10.6}$$

(b) The function $\phi_S(T)$ is a valid characteristic function for all real $\nu > 0$.

Proof. We use the approach of Muller and Chi (2006), who generalized earlier work to cover all population-singular and sample-rank-insufficient cases.

Proof of (a). For $\text{rank}(\Sigma) = p_1 = p$, Muirhead (1984) proved (10.2) for positive integer ν . The proof of Theorem 10.2 includes (10.2) as a special case ($\Omega = \mathbf{0}$) for $p_1 \leq p$. Lemma 10.1 gives (10.3). Equation (10.5) follows from $|I_p - 2iT\Phi\Phi'|^{-\nu/2} = (|\Phi^{-t} - 2iT\Phi\Phi'| |\Phi'|)^{-\nu/2} = (|\Phi'| |\Phi^{-t} - 2iT\Phi\Phi'|)^{-\nu/2}$. If $\text{rank}(\Sigma) = p = p_1$, then $\Phi' = \Phi_1'$ gives (10.4) from (10.3) and (10.6) from (10.4).

For $p_1 < p$ and $\lambda_1 > 0$, there exists $p_1 \times p_1 \ S_1 \sim \mathcal{W}_{p_1}[\nu, \text{Dg}(\lambda_1)]$ with $\phi_{S_1}(T_1) = |I_{p_1} - 2iT_1 \text{Dg}(\lambda_1)|^{-\nu/2}$ and $S = \Upsilon_1 S_1 \Upsilon_1' \sim S \mathcal{W}_p[\nu, \Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1']$. With $p \times p$ real $T = T'$ as defined in the theorem, Lemma 7.5 gives

$$\begin{aligned} \phi_S(T) &= |I_{p_1} - 2i\Upsilon_1' T (\Upsilon_1')' \text{Dg}(\lambda_1)|^{-\nu/2} \\ &= |\text{Dg}^{-1/2}(\lambda_1) - 2i\Upsilon_1' T \Upsilon_1 \text{Dg}^{1/2}(\lambda_1)|^{-\nu/2} |\text{Dg}^{1/2}(\lambda_1)|^{-\nu/2} \\ &= |I_{p_1} - 2i\text{Dg}^{1/2}(\lambda_1) \Upsilon_1' T \Upsilon_1 \text{Dg}^{1/2}(\lambda_1)|^{-\nu/2}, \end{aligned} \tag{10.7}$$

which gives equation (10.6) for $\text{rank}(\Sigma) = p_1 < p$. Equation (10.5) follows from

$$\begin{aligned}
 \phi_S(\mathbf{T}) &= (|\mathbf{I}_{p_1} - 2i\Phi_1'\mathbf{T}\Phi_1||\mathbf{I}_{p-p_1}|)^{-\nu/2} \\
 &= \left| \begin{bmatrix} \mathbf{I}_{p_1} - 2i\Phi_1'\mathbf{T}\Phi_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{p-p_1} \end{bmatrix} \right|^{-\nu/2} \\
 &= |\mathbf{I}_p - 2i[\Phi_1 \ \mathbf{0}]'\mathbf{T}[\Phi_1 \ \mathbf{0}]|^{-\nu/2} \\
 &= |\mathbf{I}_p - 2i[\Phi_1 \ \Phi_0]'\mathbf{T}[\Phi_1 \ \Phi_0]|^{-\nu/2}. \tag{10.8}
 \end{aligned}$$

With $\phi_S(\mathbf{T}) = |\mathbf{I}_{p_1} - 2i\Upsilon_1'\mathbf{T}\Upsilon_1 \text{Dg}(\lambda_1)|^{-\nu/2}$, Theorem 7.2 in Schott (2005) gives

$$\begin{aligned}
 \phi_S(\mathbf{T}) &= \left| \begin{bmatrix} \mathbf{I}_{p_1} - 2i\Upsilon_1'\mathbf{T}\Upsilon_1 \text{Dg}(\lambda_1) & \mathbf{0} \\ -2i\Upsilon_0'\mathbf{T}\Upsilon_1 \text{Dg}(\lambda_1) & \mathbf{I}_{p-p_1} \end{bmatrix} \right|^{-\nu/2} \\
 &= \left| \mathbf{I}_p - 2i \begin{bmatrix} \Upsilon_1'\mathbf{T}\Upsilon_1 & \Upsilon_1'\mathbf{T}\Upsilon_0 \\ \Upsilon_0'\mathbf{T}\Upsilon_1 & \Upsilon_0'\mathbf{T}\Upsilon_0 \end{bmatrix} \begin{bmatrix} \text{Dg}(\lambda_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \right|^{-\nu/2} \\
 &= |\mathbf{I}_p - 2i\Upsilon'\mathbf{T}\Upsilon \text{Dg}(\lambda_1, \mathbf{0})|^{-\nu/2} \\
 &= |\Upsilon\Upsilon' - 2i\Upsilon'\mathbf{T}\Upsilon \text{Dg}(\lambda_1, \mathbf{0})\Upsilon'|^{-\nu/2} \\
 &= |\mathbf{I}_p - 2i\mathbf{T}\Upsilon_1 \text{Dg}(\lambda_1)\Upsilon_1'|^{-\nu/2}. \tag{10.9}
 \end{aligned}$$

Equation (10.4) and then (10.3) and (10.2) are seen to hold for $\text{rank}(\Sigma) = p_1 < p$.

Proof of (b). Conditions guaranteeing a valid characteristic function (Kendall and Stuart, 1977, p. 105) may be verified directly for noninteger ν . \square

Theorem 10.2 (a) For $\nu > 0$, $\text{rank}(\Sigma) = \text{rank}(\Phi_1) = p_1 \leq p$, $\Sigma = \Phi_1\Phi_1'$, $p \times p_1$ Φ_1 , $i = (-1)^{1/2}$, $p \times p$ real $\mathbf{U} = \mathbf{U}'$, $\langle \mathbf{T} \rangle_{jj} = u_{jj}$, and $\langle \mathbf{T} \rangle_{jk} = u_{jk}/2$, the characteristic function of $\mathbf{S} = \mathbf{Y}'\mathbf{Y} \sim (\mathcal{S})\mathcal{W}_p(\nu, \Sigma, \Delta)$ for $\Delta = \mathbf{M}'\mathbf{M}$ and $\mathbf{Y} \sim \mathcal{N}_{\nu,p}(\mathbf{M}, \mathbf{I}_\nu, \Sigma)$ is

$$\begin{aligned}
 \phi_{S_Y}(\mathbf{T}) &= |\mathbf{I}_p - 2i\mathbf{T}\Sigma|^{-\nu/2} \exp\{i \text{tr}[\mathbf{T}\Sigma(\mathbf{I}_p - 2i\mathbf{T}\Sigma)^{-1}\Delta]\} \\
 &= |\mathbf{I}_p - 2i\mathbf{T}\Sigma|^{-\nu/2} \exp\{i \text{tr}[\mathbf{T}\Sigma(\mathbf{I}_p - 2i\mathbf{T}\Sigma)^+\Delta]\} \\
 &= |\mathbf{I}_{p_1} - 2i\Phi_1'\mathbf{T}\Phi_1|^{-\nu/2} \exp\left\{\text{tr}\left[i\mathbf{T}\Phi_1(\mathbf{I}_{p_1} - 2i\Phi_1'\mathbf{T}\Phi_1)^{-1}\Phi_1^+\Delta\right]\right\}. \tag{10.10}
 \end{aligned}$$

(b) The function $\phi_S(\mathbf{T})$ is a valid characteristic function for all real $\nu > 0$.

Proof of (a). First the result for $\Sigma = \mathbf{I}_p$ is proven. If $\mathbf{Z} \sim \mathcal{N}_{\nu,p}(\mathbf{0}, \mathbf{I}_\nu, \mathbf{I}_p)$, then $\mathbf{S}_Z = (\mathbf{Z} + \mathbf{M}_Z)'(\mathbf{Z} + \mathbf{M}_Z) \sim \mathcal{W}_p(\nu, \mathbf{I}_\nu, \mathbf{M}'_Z\mathbf{M}_Z)$. Symmetry of \mathbf{S}_Z restricts attention to $\mathbf{T} = \mathbf{V}\text{Dg}(\mathbf{t})\mathbf{V}'$ with $\mathbf{V}'\mathbf{V} = \mathbf{V}\mathbf{V}' = \mathbf{I}_p$. In turn

$$\begin{aligned}
 \phi_{S_Z}(\mathbf{T}) &= E\{\exp[\text{tr}(\mathbf{iT}'\mathbf{S}_Z)]\} \\
 &= E\{\exp\{\text{tr}[\mathbf{VDg}(\mathbf{t})\mathbf{V}'(\mathbf{Z} + \mathbf{M}_Z)'(\mathbf{Z} + \mathbf{M}_Z)]\}\} \\
 &= E\{\exp\{\mathbf{i} \text{tr}[\mathbf{Dg}(\mathbf{t})(\mathbf{ZV} + \mathbf{M}_Z\mathbf{V})'(\mathbf{ZV} + \mathbf{M}_Z\mathbf{V})]\}\} \\
 &= E\{\exp\{\mathbf{i} \text{tr}[\mathbf{Dg}(\mathbf{t})(\mathbf{U} + \mathbf{Q})'(\mathbf{U} + \mathbf{Q})]\}\} \\
 &= E\left\{\exp\left[\mathbf{i} \sum_{j=1}^p t_j (\mathbf{u}_j + \mathbf{q}_j)'(\mathbf{u}_j + \mathbf{q}_j)\right]\right\} \\
 &= E\left\{\exp\left[\mathbf{i} \sum_{j=1}^p t_j \sum_{k=1}^{\nu} (u_{kj} + q_{kj})^2\right]\right\} = E\left\{\prod_{j=1}^p \prod_{k=1}^{\nu} \exp[\mathbf{i} t_j (u_{kj} + q_{kj})^2]\right\}. \tag{10.11}
 \end{aligned}$$

Here $\mathbf{U} = \mathbf{ZV} \sim \mathcal{N}_{\nu,p}(\mathbf{0}, \mathbf{I}_{\nu}, \mathbf{I}_p)$ and $\mathbf{Q} = \mathbf{M}_Z\mathbf{V}$ is constant. Hence $\{u_{kj} + q_{kj}\}$ and $\{(u_{kj} + q_{kj})^2\}$ are independent. In turn $u_{kj} + q_{kj} \sim \mathcal{N}(q_{kj}, 1)$ implies $(u_{kj} + q_{kj})^2 \sim \chi^2(1, q_{kj}^2)$. Independence and chi-square properties give

$$\begin{aligned}
 \phi_{S_Z}(\mathbf{T}) &= \prod_{j=1}^p \prod_{k=1}^{\nu} E\{\exp[\mathbf{i}t_j(u_{kj} + q_{kj})^2]\} \\
 &= \prod_{j=1}^p \prod_{k=1}^{\nu} \left\{ (1 - 2\mathbf{i}t_j)^{-1/2} \exp[\mathbf{i}t_j q_{kj}^2 (1 - 2\mathbf{i}t_j)^{-1}] \right\} \\
 &= \left[\prod_{j=1}^p (1 - 2\mathbf{i}t_j)^{-\nu/2} \right] \exp\left[\mathbf{i} \sum_{j=1}^p t_j (1 - 2\mathbf{i}t_j)^{-1} \sum_{k=1}^{\nu} q_{kj}^2 \right] \\
 &= |\mathbf{I} - 2\mathbf{iT}|^{-\nu/2} \exp\left[\mathbf{i} \sum_{j=1}^p t_j (1 - 2\mathbf{i}t_j)^{-1} \mathbf{q}'_j \mathbf{q}_j \right] \\
 &= |\mathbf{I} - 2\mathbf{iT}|^{-\nu/2} \exp\{\mathbf{i} \text{tr}\{\mathbf{Dg}(\mathbf{t})[\mathbf{Dg}(1 - 2\mathbf{i}t_j)]^{-1} \mathbf{V}'\mathbf{M}'_Z \mathbf{M}_Z \mathbf{V}\}\} \\
 &= |\mathbf{I} - 2\mathbf{iT}|^{-\nu/2} \exp\{\mathbf{i} \text{tr}\{\mathbf{VDg}(\mathbf{t})\mathbf{V}'\mathbf{V}[\mathbf{Dg}(1 - 2\mathbf{i}t_j)]^{-1} \mathbf{V}'\mathbf{M}'_Z \mathbf{M}_Z\}\} \\
 &= |\mathbf{I} - 2\mathbf{iT}|^{-\nu/2} \exp\{\text{tr}[\mathbf{iT}(\mathbf{I} - 2\mathbf{iT})^{-1} \mathbf{M}'_Z \mathbf{M}_Z]\}. \tag{10.12}
 \end{aligned}$$

The last form is the CF of a noncentral Wishart with covariance \mathbf{I}_p .

For clarity in generalizing to $\Sigma \neq \mathbf{I}_p$ with $\text{rank}(\Sigma) = p_1 \leq p$, in the remainder of the proof \mathbf{S}_Y replaces \mathbf{S} and $\Delta = \mathbf{M}'_Y \mathbf{M}_Y$. Also $\mathbf{Z} \sim \mathcal{N}_{\nu,p_1}(\mathbf{0}, \mathbf{I}_{\nu}, \mathbf{I}_{p_1})$ has $p_1 \leq p$ columns in the remainder of the proof. For $\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{\nu,p}(\mathbf{M}_Y, \mathbf{I}_{\nu}, \Sigma)$, $p \times p_1 \Phi_1$, $\Sigma = \Phi_1 \Phi_1'$, $\nu \times p_1 \mathbf{M}_Z$, and $\mathbf{Z} \sim \mathcal{N}_{\nu,p_1}(\mathbf{0}, \mathbf{I}_{\nu}, \mathbf{I}_{p_1})$, Lemma 8.4 gives $\mathbf{Y} = (\mathbf{Z} + \mathbf{M}_Z)\Phi_1'$, with $\mathbf{M}_Y = \mathbf{M}_Z \Phi_1'$ and $\mathbf{M}_Z = \mathbf{M}_Y \Phi_1^{+t}$. In turn

$$\mathbf{S}_Y = \mathbf{Y}'\mathbf{Y} = \Phi_1(\mathbf{Z} + \mathbf{M}_Z)'(\mathbf{Z} + \mathbf{M}_Z)\Phi_1' = \Phi_1 \mathbf{S}_Z \Phi_1'. \tag{10.13}$$

If $\Phi_1 = \Upsilon_1 \mathbf{Dg}(\lambda_1)^{1/2}$, then $\Phi_1^{+t} = \Upsilon_1 \mathbf{Dg}(\lambda_1)^{-1/2}$. Lemma 7.5 gives

$$\begin{aligned}
 \phi_{S_Y}(\mathbf{T}) &= |\mathbf{I} - 2\mathbf{i}\Phi_1' \mathbf{T} \Phi_1|^{-\nu/2} \exp\left\{\text{tr}[\mathbf{i}\Phi_1' \mathbf{T} \Phi_1 (\mathbf{I}_{p_1} - 2\mathbf{i}\Phi_1' \mathbf{T} \Phi_1)^{-1} \mathbf{M}'_Z \mathbf{M}_Z]\right\} \\
 &= |\mathbf{I}_p - 2\mathbf{iT}\Sigma|^{-\nu/2} \exp\left\{\text{tr}[\mathbf{iT}\Phi_1 (\mathbf{I}_{p_1} - 2\mathbf{i}\Phi_1' \mathbf{T} \Phi_1)^{-1} \mathbf{M}'_Z \mathbf{M}_Y]\right\}. \tag{10.14}
 \end{aligned}$$

The following equalities arise from repeated use of (1) $\Phi_1 = \Upsilon_1 \text{Dg}(\lambda_1)^{1/2}$ in Lemma 10.1, (2) $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$ for square and full-rank matrices, and (3) similar generalized inverse results in parts (a) and (b) of Lemma 1.15.

$$\begin{aligned}
 & T\Phi_1 \left[\text{Dg}(\lambda_1)^{-1/2} - 2i\Upsilon_1' T\Phi_1 \right]^{-1} \text{Dg}(\lambda_1)^{-1/2} \mathbf{M}'_Z \mathbf{M}_Y = \\
 & T\Phi_1 \left[\Upsilon_1' \Upsilon_1 \text{Dg}(\lambda_1)^{-1/2} - 2i\Upsilon_1' T\Phi_1 \right]^{-1} \text{Dg}(\lambda_1)^{-1/2} \mathbf{M}'_Z \mathbf{M}_Y = \\
 & T\Phi_1 \left[\Upsilon_1 \text{Dg}(\lambda_1)^{-1/2} - 2iT\Phi_1 \right]^+ \Upsilon_1' \text{Dg}(\lambda_1)^{-1/2} \mathbf{M}'_Z \mathbf{M}_Y = \\
 & T\Phi_1 \left\{ [\Upsilon_1 - 2iT\Upsilon_1 \text{Dg}(\lambda_1)] \text{Dg}(\lambda_1)^{-1/2} \right\}^+ \mathbf{M}'_Y \mathbf{M}_Y = \\
 & T\Upsilon_1 \text{Dg}(\lambda_1) \{ [\Upsilon_1 - 2iT\Upsilon_1 \text{Dg}(\lambda_1)] \}^+ \Delta = \\
 & T\Upsilon_1 \text{Dg}(\lambda_1) [\Upsilon_1 - 2iT\Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1' \Upsilon_1]^+ \Delta = \\
 & T\Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1' [I_p - 2iT\Upsilon_1 \text{Dg}(\lambda_1) \Upsilon_1']^+ \Delta = T\Sigma (I_p - 2iT\Sigma)^+ \Delta. \quad (10.15)
 \end{aligned}$$

For $p \times p$ Υ , $(I_p - 2iT\Sigma)$ is similar to $\mathbf{A} = [I_p - 2i\Upsilon' T \Upsilon \text{Dg}(\lambda_1, \mathbf{0})]$. For $\mathbf{D} = \text{Dg}(\lambda_1, \mathbf{1})^{1/2}$ and T partitioned like Υ , $(I_p - 2iT\Sigma)$ and \mathbf{A} are similar to

$$\mathbf{DAD}^{-1} = \begin{bmatrix} I_{p_1} - 2i\Phi_1' T_{1,1} \Phi_1 & \mathbf{0} \\ T_{2,1} \Phi_1 & I_{p-p_1} \end{bmatrix}. \quad (10.16)$$

The triangular form ensures $\text{rank}(\mathbf{DAD}^{-1}) = \text{rank}(I_{p_1} - 2i\Phi_1' T_{1,1} \Phi_1) + (p - p_1)$. In turn, $\Phi_1' T_{1,1} \Phi_1 = \mathbf{V} \text{Dg}(\mathbf{a}) \mathbf{V}'$ and $(I_{p_1} - 2i\Phi_1' T_{1,1} \Phi_1) = \mathbf{V} [I_{p_1} - 2i \text{Dg}(\mathbf{a})] \mathbf{V}'$. Consequently $\text{rank}(I_p - 2iT\Sigma) < p$ only if $(1 - 2ia_k) = 0$ for one or more a_k , which never happens because a_k is always real. Furthermore $(I_p - 2iT\Sigma)^+ = (I_p - 2iT\Sigma)^{-1}$. A parallel analysis for the moment generating function gives $(1 - 2a_k) > 0$ as sufficient to ensure existence and full-rank $(I_p - 2iT\Sigma)$.

Proof of (b). Conditions guaranteeing a function is a valid characteristic function (Kendall and Stuart, 1977, p. 105) may be verified for noninteger ν . \square

10.3 PROPERTIES OF THE WISHART

Theorem 10.3 For $\nu \times p$ $\mathbf{Y} = [\mathbf{y}_1 \mathbf{y}_2 \cdots \mathbf{y}_p] \sim \mathcal{N}_{\nu,p}(\mathbf{0}, I_\nu, \Sigma)$, $\text{rank}(\Sigma) = p$, $\nu \geq p$, the $p \times p$ matrix $\mathbf{S} = \mathbf{Y}'\mathbf{Y}$ has a central nonsingular Wishart distribution with parameters ν and Σ . The joint PDF of the $p(p+1)/2$ distinct elements of \mathbf{S} is (Gupta and Nagar, 2000)

$$f_{\mathbf{S}}(\mathbf{S}_*) = \frac{|\mathbf{S}_*|^{(\nu-p-1)/2} \exp[-\text{tr}(\Sigma^{-1} \mathbf{S}_*)/2]}{2^{\nu p/2} \pi^{p(p-1)/4} |\Sigma|^{\nu/2} \prod_{j=1}^p \Gamma[(\nu+1-j)/2]}. \quad (10.17)$$

The j, j' element of \mathbf{S} is a bilinear form for $j \neq j'$ and a univariate quadratic form if $j = j'$. If $\nu < p$ the density does not exist.

With $\nu = N - \text{rank}(\mathbf{X})$, the central Wishart describes the distribution of $\nu \widehat{\Sigma}$ in the multivariate general linear model. The noncentral form applies to the hypothesis sum-of-squares matrix under the alternative.

Theorem 10.4 If $\mathbf{S} \sim (S)\mathcal{W}_p(\nu, \Sigma, \Delta)$ and \mathbf{T} is any $p \times p_1$ constant, then

$$\mathbf{T}'\mathbf{S}\mathbf{T} \sim (S)\mathcal{W}_{p_1}(\nu, \mathbf{T}'\Sigma\mathbf{T}, \mathbf{T}'\Delta\mathbf{T}). \quad (10.18)$$

Proof. Left as an exercise. *Hints:* characteristic function, linear transformation.

Many authors consider only full-rank Σ and \mathbf{T} (and $p_1 \leq p$). The result holds for any sort of Wishart and any conforming constant \mathbf{T} . The transformation is equivalent to having transformed the underlying Gaussian (which can be the basis of a proof). If $\mathbf{S} = \mathbf{Y}'\mathbf{Y}$ with $\mathbf{Y} \sim (S)\mathcal{N}_{\nu,p}(\mathbf{M}, \mathbf{I}_\nu, \Sigma)$, then

$$\mathbf{T}'\mathbf{S}\mathbf{T} = \mathbf{T}'\mathbf{Y}'\mathbf{Y}\mathbf{T} = \mathbf{Y}'_1\mathbf{Y}_1 \quad (10.19)$$

and

$$\mathbf{Y}_1 = \mathbf{Y}\mathbf{T} \sim \mathcal{N}_{\nu,p_1}(\mathbf{M}\mathbf{T}, \mathbf{I}_\nu, \mathbf{T}'\Sigma\mathbf{T}). \quad (10.20)$$

Theorem 10.5 If $\mathbf{S} \sim \mathcal{W}_p(\nu, \Sigma, \Delta)$ and \mathbf{t} is a vector of constants, then $\mathbf{t}'\mathbf{S}\mathbf{t}/\mathbf{t}'\Sigma\mathbf{t} \sim \chi^2(\nu, \omega)$, with $\omega = \mathbf{t}'\Delta\mathbf{t}/\mathbf{t}'\Sigma\mathbf{t}$.

Proof. With $\mathbf{S} = \sum_{i=1}^{\nu} \mathbf{Y}_i'\mathbf{Y}_i$ for $[\text{row}_i(\mathbf{Y})]' = \mathbf{Y}_i' \sim \mathcal{N}_p(\boldsymbol{\mu}_i, \Sigma)$ it follows that

$$\mathbf{t}'\mathbf{S}\mathbf{t} = \mathbf{t}'\left(\sum_{i=1}^{\nu} \mathbf{Y}_i'\mathbf{Y}_i\right)\mathbf{t} = \sum_{i=1}^{\nu} \mathbf{t}'\mathbf{Y}_i'\mathbf{Y}_i\mathbf{t} = \sum_{i=1}^{\nu} (\mathbf{t}'\mathbf{Y}_i')^2. \quad (10.21)$$

Also $\mathbf{t}'\mathbf{Y}_i \sim \mathcal{N}_1(\mathbf{t}'\boldsymbol{\mu}_i, \mathbf{t}'\Sigma\mathbf{t})$ allows concluding

$$\mathbf{t}'\mathbf{Y}_i/(\mathbf{t}'\Sigma\mathbf{t})^{1/2} \sim \mathcal{N}_1[\mathbf{t}'\boldsymbol{\mu}_i/(\mathbf{t}'\Sigma\mathbf{t})^{1/2}, 1]. \quad (10.22)$$

The result follows because a squared unit Gaussian is a noncentral chi square. \square

The result can be proven more generally as a special case of the previous theorem. The converse is *not* true (Mitra, 1970). Theorem 10.9 provides a form of a converse, based on a far stronger condition.

Theorem 10.6 Principal diagonal blocks of a Wishart are Wishart and independent if and only if the corresponding interblock covariance is zero. If

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix} \sim (S)\mathcal{W}_p(\nu, \Sigma, \Delta), \quad (10.23)$$

with Σ and Δ partitioned to match, \mathbf{S}_{11} ($q \times q$), then

$$\mathbf{S}_{11} \sim (\mathcal{S})\mathcal{W}_q(\nu, \boldsymbol{\Sigma}_{11}, \boldsymbol{\Delta}_{11}) \tag{10.24}$$

$$\mathbf{S}_{22} \sim (\mathcal{S})\mathcal{W}_{p-q}(\nu, \boldsymbol{\Sigma}_{22}, \boldsymbol{\Delta}_{22}) \tag{10.25}$$

$$\mathbf{S}_{11} \perp\!\!\!\perp \mathbf{S}_{22} \Leftrightarrow \boldsymbol{\Sigma}_{12} = \mathbf{0}. \tag{10.26}$$

Proof. Highly recommended as an exercise.

Theorem 10.7 If $\mathbf{S}_1 \sim (\mathcal{S})\mathcal{W}_p(\nu_1, \boldsymbol{\Sigma}, \boldsymbol{\Delta}_1) \perp\!\!\!\perp \mathbf{S}_2 \sim (\mathcal{S})\mathcal{W}_p(\nu_2, \boldsymbol{\Sigma}, \boldsymbol{\Delta}_2)$ then $\mathbf{S}_1 + \mathbf{S}_2 \sim (\mathcal{S})\mathcal{W}_p(\nu_1 + \nu_2, \boldsymbol{\Sigma}, \boldsymbol{\Delta}_1 + \boldsymbol{\Delta}_2)$.

Proof. Highly recommended as an exercise.

Theorem 10.8 If $\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \boldsymbol{\Sigma})$, $\text{rank}(\boldsymbol{\Sigma}) = p_1 \leq p$ with $N \geq p_1$, $N \times N$ \mathbf{A} and \mathbf{B} are constants, the following hold.

- (a) \mathbf{AY} and \mathbf{BY} are independent if and only if $\mathbf{AB}' = \mathbf{0}$.
- (b) If $\mathbf{A} = \mathbf{A}'$ is positive definite or positive semidefinite and $\mathbf{AB}' = \mathbf{0}$, then \mathbf{BY} and $\mathbf{Y}'\mathbf{AY}$ are independent.
- (c) If $\mathbf{A} = \mathbf{A}'$ and $\mathbf{B} = \mathbf{B}'$ are positive definite or positive semidefinite and $\mathbf{BA} = \mathbf{0}$, then $\mathbf{Y}'\mathbf{AY}$ and $\mathbf{Y}'\mathbf{BY}$ are independent.
- (d) If $\mathbf{A} = \mathbf{A}' = \mathbf{A}^2$, then $\nu = \text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A})$ and

$$\mathbf{S} = \mathbf{Y}'\mathbf{AY} \sim (\mathcal{S})\mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \mathbf{M}'\mathbf{A}\mathbf{M}). \tag{10.27}$$

Proof. The basic approach is as follows. Here $\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \boldsymbol{\Sigma})$ gives

$$\begin{bmatrix} \mathbf{Y} \\ \mathbf{Y} \end{bmatrix} \sim (\mathcal{S})\mathcal{N}_{2N,p} \left(\begin{bmatrix} \mathbf{M} \\ \mathbf{M} \end{bmatrix}, [(\mathbf{1}_2\mathbf{1}'_2) \otimes \mathbf{I}_N], \boldsymbol{\Sigma} \right) \tag{10.28}$$

$$\begin{bmatrix} \mathbf{AY} \\ \mathbf{BY} \end{bmatrix} \sim (\mathcal{S})\mathcal{N}_{(\nu_a + \nu_b),p} \left(\begin{bmatrix} \mathbf{M} \\ \mathbf{M} \end{bmatrix}, \begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} [(\mathbf{1}_2\mathbf{1}'_2) \otimes \mathbf{I}_N] \begin{bmatrix} \mathbf{A}' & \mathbf{0} \\ \mathbf{0} & \mathbf{B}' \end{bmatrix}, \boldsymbol{\Sigma} \right) \tag{10.29}$$

$$\begin{bmatrix} \mathbf{A} & \mathbf{0} \\ \mathbf{0} & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{I}_N & \mathbf{I}_N \\ \mathbf{I}_N & \mathbf{I}_N \end{bmatrix} \begin{bmatrix} \mathbf{A}' & \mathbf{0} \\ \mathbf{0} & \mathbf{B}' \end{bmatrix} = \begin{bmatrix} \mathbf{AA}' & \mathbf{AB}' \\ \mathbf{BA}' & \mathbf{BB}' \end{bmatrix}. \tag{10.30}$$

Part (a) follows immediately. If $\mathbf{A} = \mathbf{A}'$ is nonnegative definite, then $\mathbf{A} = \mathbf{F}\mathbf{F}'$ and $\mathbf{Y}'\mathbf{AY} = \mathbf{Y}'\mathbf{F}\mathbf{F}'\mathbf{Y}$. Parts (b) and (c) follow from independence of the underlying Gaussian matrices, as seen from part (a). In part (d), $\mathbf{A} = \mathbf{A}' = \mathbf{A}^2$ implies $\mathbf{A} = \mathbf{V}\mathbf{V}'$ with $\mathbf{V}'\mathbf{V} = \mathbf{I}_\nu$, $\mathbf{Y}_A = \mathbf{V}'\mathbf{Y} \sim (\mathcal{S})\mathcal{N}_{N,\nu}(\mathbf{V}'\mathbf{M}, \mathbf{I}_\nu, \boldsymbol{\Sigma})$, and $\mathbf{S} = \mathbf{Y}'_A\mathbf{Y}_A \sim (\mathcal{S})\mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \mathbf{M}'\mathbf{V}\mathbf{V}'\mathbf{M})$. \square

What does the theorem tell us about the independence of the usual sample covariance and sample mean vector?

Corollary 10.8.1 If $p = 1$ and $\sigma^2 > 0$, then $\mathbf{y} \sim \mathcal{N}_{n,1}(\boldsymbol{\mu}, \mathbf{I}_n, \sigma^2) \Leftrightarrow \mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \mathbf{I}_n \sigma^2)$ and

$$s = \mathbf{y}' \mathbf{A} \mathbf{y} \sim \mathcal{W}_1(\nu, \sigma^2, \boldsymbol{\mu}' \mathbf{A} \boldsymbol{\mu}). \quad (10.31)$$

Equivalently

$$s/\sigma^2 \sim \chi^2(\nu, \boldsymbol{\mu}' \mathbf{A} \boldsymbol{\mu}/\sigma^2). \quad (10.32)$$

Proof. Left as an exercise.

Corollary 10.8.2 The matrix $\mathbf{Y}' \mathbf{A} \mathbf{Y}$ is Wishart if and only if $\mathbf{A} = \mathbf{A}' = \mathbf{A}^2$ (idempotent).

Proof. Only the statement “ $\mathbf{Y}' \mathbf{A} \mathbf{Y}$ is Wishart implies \mathbf{A} is idempotent” must be proven. Spectral decomposition gives $\mathbf{A} = \mathbf{V}_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{V}_1'$ with \mathbf{V}_1 $N \times \nu$ and $\mathbf{V}_1' \mathbf{V}_1 = \mathbf{I}_\nu$. In turn $\mathbf{Y}_1 = \mathbf{V}_1' \mathbf{Y} \sim \mathcal{N}_{\nu,p}(\mathbf{V}_1' \mathbf{M}, \mathbf{I}_\nu, \boldsymbol{\Sigma})$ and

$$\mathbf{Y}' \mathbf{A} \mathbf{Y} = \mathbf{Y}' \mathbf{V}_1 \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{V}_1' \mathbf{Y} = \mathbf{Y}_1' \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{Y}_1. \quad (10.33)$$

The last equation can only be guaranteed to exist if $\mathbf{A} = \mathbf{A}'$. By the definition of a Wishart, $\mathbf{Y}_1' \mathbf{Y}_1 \sim \mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \mathbf{M}' \mathbf{V}_1 \mathbf{V}_1' \mathbf{M})$. Referring again to the definition of a Wishart, the $\nu \times \nu$ matrix $\text{Dg}(\boldsymbol{\lambda}_1)$ must equal \mathbf{I}_ν for $\mathbf{Y}_1' \text{Dg}(\boldsymbol{\lambda}_1) \mathbf{Y}_1 = \mathbf{Y}' \mathbf{A} \mathbf{Y}$ to be a Wishart. Requiring $\mathbf{A} = \mathbf{V}_1 \mathbf{I}_\nu \mathbf{V}_1'$ is equivalent to requiring $\mathbf{A} = \mathbf{A}' = \mathbf{A}^2$. \square

If all conditions for part (a) of the last theorem are met except $\mathbf{A} \neq \mathbf{A}^2$, then $\mathbf{Y}' \mathbf{A} \mathbf{Y}$ is a general multivariate quadratic form. As a generalization of results for univariate quadratic forms, a constituent matrix decomposition of \mathbf{A} gives

$$\mathbf{Y}' \mathbf{A} \mathbf{Y} = \sum_{k=1}^{\nu} \lambda_{1k} \mathbf{Y}' \mathbf{v}_{1k} \mathbf{v}_{1k}' \mathbf{Y} \quad (10.34)$$

with i.i.d. $\mathbf{Y}' \mathbf{v}_{1k} \mathbf{v}_{1k}' \mathbf{Y} \sim (\mathcal{S}) \mathcal{W}_p(1, \boldsymbol{\Sigma}, \mathbf{M} \mathbf{v}_{1k} \mathbf{v}_{1k}' \mathbf{M})$. We leave consideration of such forms for another venue.

Theorem 10.9 For $N \times p$ $\mathbf{Y} = [\mathbf{y}_1 \mathbf{y}_2 \cdots \mathbf{y}_p] \sim \mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \boldsymbol{\Sigma})$, $\text{rank}(\boldsymbol{\Sigma}) = p$, $N \geq p$, and $\mathbf{Y}_i = \text{row}_i(\mathbf{Y})$. The constant $N \times N$ matrix $\mathbf{A} = \mathbf{A}'$ has $\nu = \text{rank}(\mathbf{A})$ and $\boldsymbol{\Delta} = \mathbf{M}' \mathbf{A} \mathbf{M}$. For the conditions given,

$$\mathbf{S} = \mathbf{Y}' \mathbf{A} \mathbf{Y} \sim \mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \boldsymbol{\Delta}) \quad (10.35)$$

if and only if for all $\mathbf{t} \in \mathbb{R}^p$

$$\frac{\mathbf{t}' \mathbf{Y}' \mathbf{A} \mathbf{Y} \mathbf{t}}{\mathbf{t}' \boldsymbol{\Sigma} \mathbf{t}} \sim \chi^2(\nu, \mathbf{t}' \boldsymbol{\Delta} \mathbf{t} / \mathbf{t}' \boldsymbol{\Sigma} \mathbf{t}). \quad (10.36)$$

Furthermore $\nu = \text{tr}(\mathbf{A})$.

Proof. The fact $t'Y'AYt/t'\Sigma t$ is chi square may be proven by combining two previous theorems. The proof for the other direction begins with defining $N \times 1$ $x = Yt$ for a fixed nonzero t and $Yt \sim \mathcal{N}_{N,1}(Mt, I_N, t'\Sigma t)$. Equivalently $x \sim \mathcal{N}_N(Mt, I_N t'\Sigma t)$. By assumption $t'Y'AYt/(t'\Sigma t) = x'Ax/\mathcal{V}(x) \sim \chi^2(\nu, t'\Delta t/t'\Sigma t)$. By Theorem 9.6, for $x'Ax/\mathcal{V}(x) \sim \chi^2(\nu, t'\Delta t/t'\Sigma t)$ to hold, the eigenvalues of A must all be 1. Therefore A is idempotent, which suffices, with Theorem 10.8, to ensure $S \sim \mathcal{W}_p(\nu, \Sigma, \Delta)$. \square

Lemma 10.2 (a) If $S \sim (S)\mathcal{W}_p(\nu, \Sigma, \Delta)$, then, without loss of generality, it may be assumed $S = Y_1'Y_1$ with $Y_1 \sim (S)\mathcal{N}_{\nu,p}(M, I_\nu, \Sigma)$.

(b) If $A = A' = A^2$, then $A = V_1V_1'$ with V_1 $N \times \nu$ and $V_1'V_1 = I_\nu$.

If $Y \sim (S)\mathcal{N}_{N,p}(M, I_N, \Sigma)$, then $Y_1 = V_1'Y \sim (S)\mathcal{N}_{\nu,p}(M, I_\nu, \Sigma)$ and $Y'AY = Y_1'Y_1 \sim (S)\mathcal{W}_p(\nu, \Sigma, M'V_1V_1'M)$.

Proof. (a) By definition of a Wishart. **(b)** Shown by construction. \square

Theorem 10.10 For $N \times p$ $Y = [y_1 \ y_2 \ \cdots \ y_p] \sim \mathcal{N}_{N,p}(M, I_N, \Sigma)$, $N \geq p$, $\text{rank}(\Sigma) = p$, and $Y_i = \text{row}_i(Y)$. The constant $N \times N$ matrix $A = A'$ has $\nu = \text{rank}(A)$ and $\Delta = M'AM$. Given the definitions,

$$E(Y'AY) = \text{tr}(A)\Sigma + \Delta. \tag{10.37}$$

The result remains true with the weaker assumption of i.i.d. rows with finite second moments, even without Gaussian variables.

The ANOVA theorem generalizes directly from univariate to multivariate quadratic forms. The result is the MANOVA theorem.

Theorem 10.11 (MANOVA theorem) The $N \times p$ matrix $Y = [y_1 \ y_2 \ \cdots \ y_p]$, with $Y_i = \text{row}_i(Y)$ and $N \geq p$, is matrix Gaussian, $Y \sim \mathcal{N}_{N,p}(M, I_N, \Sigma)$, with $\text{rank}(\Sigma) = p$. Also, with $k \in \{1, 2, \dots, d\}$ and $N \times N$ matrix $A_k = A'_k$ of rank $\nu_k > 0$, the matrix $A_0 = \sum_{k=1}^d A_k$, with rank of ν_0 , is symmetric and $N \times N$. The assumptions allow defining five conditions.

1. A_k is idempotent for $k \in \{1, 2, \dots, d\}$;
2. $A_k A_{k'} = 0 \ \forall k \neq k'$ with $k \in \{1, 2, \dots, d\}$ and $k' \in \{1, 2, \dots, d\}$;
3. A_0 is idempotent;
4. $\nu_0 = \sum_{k=1}^d \nu_k$, which is equivalent to $\text{rank}(\sum_{k=1}^d A_k) = \sum_{k=1}^d \text{rank}(A_k)$; and
5. A_k is idempotent for $k \in \{1, \dots, d-1\}$ and A_d is nonnegative definite.

Given the assumptions,

(a) $Y'A_k Y \sim \mathcal{W}_p(\nu_k, \Sigma, M'A_k M)$ for $k \in \{0, 1, 2, \dots, d\}$ and

(b) $\{Y'A_k Y\}$ are mutually independent for $k \in \{1, 2, \dots, d\}$ and $k \neq 0$ if and only if

- I. any two of conditions 1, 2 and 3 are true, or
- II. conditions 3 and 4 are true, or

III. conditions 3 and 5 are true.
 Furthermore, if *any* of I, II, or III, holds, then *all* hold.

Proof. Mostly left as an exercise. Theorem 10.9 combines with the ANOVA theorem and Theorem 5 in Searle (1971, Section 2.5).

Theorem 10.12 The conditional sum-of-squares matrix from a Wishart is Wishart. If

$$\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix} \sim \mathcal{W}_p(\nu, \Sigma), \tag{10.38}$$

with Σ partitioned to match, full-rank \mathbf{S}_{22} , and $\mathbf{S}_{1.2} = \mathbf{S}_{11} - \mathbf{S}_{12}\mathbf{S}_{22}^{-1}\mathbf{S}_{21}$, then

$$\mathbf{S}_{1.2} \sim \mathcal{W}_q(\nu - p + q, \Sigma_{1.2}). \tag{10.39}$$

Proof. Highly recommended as an exercise.

Corollary 10.12 For any combination of singular Σ_{22} and $\nu \leq p$, it follows that $\mathbf{S}_{1.2} = \mathbf{S}_{11} - \mathbf{S}_{12}\mathbf{S}_{22}^- \mathbf{S}_{21} = \mathbf{S}_{11} - \mathbf{S}_{12}\mathbf{S}_{22}^+ \mathbf{S}_{21}$ and $\mathbf{S}_{1.2} \sim (\mathcal{S})\mathcal{W}_q(\nu - p + q, \Sigma_{1.2})$.

Theorem 10.13 (Wijsman, 1959; Odell and Feiveson, 1966) If $\mathbf{S} \sim \mathcal{W}_p(\nu, \mathbf{I}, \mathbf{0})$ and $\nu \geq p$, then a $p \times p$ lower triangular matrix \mathbf{T} exists such that $\mathbf{S} = \mathbf{T}\mathbf{T}'$ and

1. $t_{jj} \sim \chi(\nu - j + 1, 0)$ for the diagonal elements,
2. $t_{jj'} \sim \mathcal{N}_1(0, 1)$ for $j > j'$, the lower off-diagonal elements, and
3. all elements of \mathbf{T} are statistically independent.

The result, the *Bartlett decomposition*, also holds for any (finite) real $\nu > p$.

Corollary 10.13 (a) If $\Sigma = \Phi\Phi'$ is positive definite, \mathbf{T}_1 has the properties described in the theorem, $\mathbf{S}_1 = \mathbf{T}_1\mathbf{T}'_1 \sim \mathcal{W}_p(\nu, \mathbf{I})$, $\nu \geq p$, and $\mathbf{T}_2 = \Phi\mathbf{T}_1$, then $\mathbf{S}_2 = \mathbf{T}_2\mathbf{T}'_2 \sim \mathcal{W}_p(\nu, \Sigma)$.

(b) If $\nu \geq p + 1$ and $\mathbf{y} \sim \mathcal{N}_p(\boldsymbol{\mu}, \Sigma)$, then a noncentral Wishart, $\mathbf{S}_+ \sim \mathcal{W}_p(\nu, \Sigma, \boldsymbol{\mu}\boldsymbol{\mu}')$, may be generated as $\mathbf{S}_+ = \mathbf{S}_1 + \mathbf{S}_2$, with $\mathbf{S}_1 \sim \mathcal{W}_p(\nu - 1, \Sigma, \mathbf{0})$ and $\mathbf{S}_2 = \mathbf{y}\mathbf{y}' \sim \mathcal{W}_p(1, \Sigma, \boldsymbol{\mu}\boldsymbol{\mu}')$. Furthermore, as long as $\nu > p + 1$, the approach also allows directly generating noncentral pseudo-random Wishart matrices with fractional degrees of freedom.

The corollary leads to a simulation algorithm for generating pseudo-random Wishart matrices that can be dramatically faster than computing the inner product of a matrix of pseudo-random Gaussian variables. As long as $\nu > p$, the approach also allows generating pseudo-random Wishart matrices with fractional ν .

The steps are as follows.

1. Compute Φ via a Cholesky algorithm or a spectral decomposition. If $\Sigma = \Upsilon\text{Dg}(\lambda)\Upsilon'$ then $\Phi = \Upsilon\text{Dg}(\lambda)^{1/2}$.
2. Generate the required $p(p + 1)/2$ nonzero elements for matrix \mathbf{T}_1 using pseudo-

random number generators.

3. Compute $T_2 = \Phi T_1$, then $T_2 T_2'$, a pseudo-random realization.

Repeat steps 2 and 3 to create the desired number of replications. For fixed p , as ν becomes large, the approach is faster than generating a $\nu \times p$ matrix of pseudo-random Gaussian variables and taking the inner product.

10.4 THE INVERSE WISHART

The inverse of a Wishart matrix occurs in many forms of multivariate test statistics. As a special case, the F statistic for testing a general linear hypothesis in the univariate GLM may be written, with $r = \text{rank}(\mathbf{X})$,

$$F(\mathbf{y}) = \frac{SSH/a}{SSE/(N-r)} = SSH(SSE)^{-1}[(N-r)/a]. \tag{10.40}$$

Here $SSE = \sigma^2 x_e \sim \mathcal{W}_1(N-r, \sigma^2, 0)$ and $x_e \sim \chi^2(N-r, 0)$.

Definition 10.3 For $\nu \geq p = \text{rank}(\Sigma)$, the $p \times p$ matrix \mathbf{T} has the inverse (central) Wishart distribution, denoted $\mathbf{T} \sim \mathcal{W}_p^{-1}(\nu, \Sigma^{-1})$, if and only if $\mathbf{T}^{-1} = \mathbf{S} \sim \mathcal{W}_p(\nu, \Sigma)$.

Certain features of the definition should be noted. Most importantly, the matrices \mathbf{T} and Σ are implicitly assumed to be symmetric, full rank, and positive definite. Furthermore, the definition leaves open the possibility of extending the concept to noncentral Wishart matrices. Finally, the random matrix is defined in terms of another random matrix, with no reference to the distribution function or other properties. Given the conditions of the definition, \mathbf{T} has the density

$$f(\mathbf{T}; \nu, \Sigma^{-1}) = \frac{|\mathbf{T}|^{-(\nu+p+1)/2} \exp[-\text{tr}(\Sigma^{-1}\mathbf{T}^{-1})/2]}{2^{\nu p/2} \pi^{p(p-1)/4} |\Sigma|^{\nu/2} \prod_{j=1}^p \Gamma[(\nu+1-j)/2]}. \tag{10.41}$$

Mardia, Kent, and Bibby (1979, p. 85) gave a derivation. It is simple to prove $E(\mathbf{T}) = [\Sigma(\nu-p-1)]^{-1}$. Sampson (1974) stated the following two lemmas.

Lemma 10.3 If $\mathbf{T} \sim \mathcal{W}_p^{-1}(\nu, \Sigma^{-1})$ and constant \mathbf{A} is full rank and $p \times p$, then $\mathbf{A}'\mathbf{T}\mathbf{A} \sim \mathcal{W}_p^{-1}(\nu, \mathbf{A}'\Sigma^{-1}\mathbf{A})$.

Proof. Left as an exercise.

The result can be generalized, at least to $p \times p_1$ \mathbf{A} of rank p_1 with $p_1 < p$.

Lemma 10.4 If \mathbf{T}_{jk} and Σ_{jk} are $p_j \times p_k$ with $p_j + p_k = p$ and also

$$\mathbf{T} = \begin{bmatrix} \mathbf{T}_{11} & \mathbf{T}_{12} \\ \mathbf{T}_{21} & \mathbf{T}_{22} \end{bmatrix} \sim \mathcal{W}_p^{-1} \left(\nu, \begin{bmatrix} \Psi_{11} & \Psi_{12} \\ \Psi_{21} & \Psi_{22} \end{bmatrix} \right), \tag{10.42}$$

then $\mathbf{T}_{11} \sim \mathcal{W}_{p_1}^{-1}(\nu - p_2, \Psi_{11})$.

Proof. Left as an exercise. *Hint:* Use theorems about the inverse of a partitioned matrix and $\mathbf{S}_{1.2} = \mathbf{S}_{11} - \mathbf{S}_{12}\mathbf{S}_{22}^{-1}\mathbf{S}_{21}$.

10.5 RELATED DISTRIBUTIONS

The Wishart has useful relationships to other distributions. We have earlier noted $\mathbf{S} \sim \mathcal{W}_p(\nu, \Sigma, \Delta) \Leftrightarrow q_{\mathbf{t}} = \mathbf{t}'\mathbf{S}\mathbf{t}/(\mathbf{t}'\Sigma\mathbf{t}) \sim \chi^2[\nu, \mathbf{t}'\Delta\mathbf{t}/(\mathbf{t}'\Sigma\mathbf{t})] \quad \forall \mathbf{t} \in \Re^p$. The trace and determinant of a Wishart have simple distributions.

Theorem 10.14 Idempotent matrix \mathbf{A} of rank $\nu > 0$ is constant and $N \times N$. Also $p \times p$ Σ is symmetric, positive semidefinite, rank $p_1 \leq p$, with spectral decomposition $\Sigma = \mathbf{T}\text{Dg}(\lambda)\mathbf{T}'$, and $\mathbf{T}'\mathbf{T} = \mathbf{I}_p$. If λ_1 is the $p_1 \times 1$ vector of strictly positive eigenvalues, then, without loss of generality,

$$\Sigma = [\mathbf{T}_1 \ \mathbf{T}_0]\text{Dg}(\lambda_1, \mathbf{0}) \begin{bmatrix} \mathbf{T}'_1 \\ \mathbf{T}'_0 \end{bmatrix}. \tag{10.43}$$

We indicate a $p \times 1$ vector with a 1 in row k and 0 elsewhere as \mathbf{d}_k . If $\mathbf{Y} = (\mathcal{S})\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \Sigma)$, $\mathbf{m}_k = \mathbf{M}\mathbf{T}\mathbf{d}_k$ and $x_k \sim \chi^2(\nu, \mathbf{m}'_k\mathbf{A}\mathbf{m}_k/\lambda_k)$ are independent, then

$$\begin{aligned} \text{tr}(\mathbf{Y}'\mathbf{A}\mathbf{Y}) &= \sum_{k=1}^{p_1} \lambda_k x_k + \sum_{k=p_1+1}^p \mathbf{m}'_k\mathbf{A}\mathbf{m}_k \\ &= \sum_{k=1}^{p_1} \lambda_k x_k + \text{tr}(\mathbf{T}'_0\mathbf{M}'\mathbf{A}\mathbf{M}\mathbf{T}_0). \end{aligned} \tag{10.44}$$

If $p_1 = p$, then $\text{tr}(\mathbf{T}'_0\mathbf{M}'\mathbf{M}\mathbf{T}_0) = 0$. If $\mathbf{M} = \mathbf{0}$, then $\text{tr}(\mathbf{T}'_0\mathbf{M}'\mathbf{M}\mathbf{T}_0) = 0$ and $x_k \sim \chi^2(\nu, 0)$.

Proof. Glueck and Muller (1998) proved the result.

Corollary 10.14.1 The reproductive property of the matrix Gaussian allow concluding

$$\begin{aligned} [\mathbf{Y}_1 \ \mathbf{Y}_0] &= \mathbf{Y}[\mathbf{T}_1 \ \mathbf{T}_0] \\ &= \mathbf{Y}_1\mathbf{T}'_1 + \mathbf{Y}_0\mathbf{T}'_0 \sim (\mathcal{S})\mathcal{N}_{N,p}[\mathbf{M}\mathbf{T}, \mathbf{I}_N, \text{Dg}(\lambda_1, \mathbf{0})]. \end{aligned} \tag{10.45}$$

Therefore \mathbf{Y}_1 and \mathbf{Y}_0 are statistically independent with

$$\mathbf{Y}_1 \sim \mathcal{N}_{N,p_1}[\mathbf{M}\boldsymbol{\Upsilon}_1, \mathbf{I}_N, \text{Dg}(\boldsymbol{\lambda}_1)] \quad (10.46)$$

$$\mathbf{Y}_0 \sim \mathcal{SN}_{N,p-p_1}[\mathbf{M}\boldsymbol{\Upsilon}_0, \mathbf{I}_N, \text{Dg}(\mathbf{0})]. \quad (10.47)$$

In turn

$$\mathbf{Y}'_1 \mathbf{A} \mathbf{Y}_1 \sim \mathcal{W}_{p_1}[\nu, \text{Dg}(\boldsymbol{\lambda}_1), \boldsymbol{\Upsilon}'_1 \mathbf{M}' \mathbf{A} \mathbf{M} \boldsymbol{\Upsilon}_1] \quad (10.48)$$

$$\mathbf{Y}'_0 \mathbf{A} \mathbf{Y}_0 \sim \mathcal{W}_{p-p_1}[\nu, \text{Dg}(\mathbf{0}), \boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{A} \mathbf{M} \boldsymbol{\Upsilon}_0], \quad (10.49)$$

with $\mathbf{Y}'_1 \mathbf{A} \mathbf{Y}_1$ independent of (degenerate and discrete) $\mathbf{Y}'_0 \mathbf{A} \mathbf{Y}_0$. In fact, $\Pr\{\mathbf{Y}'_0 \mathbf{A} \mathbf{Y}_0 = \boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{A} \mathbf{M} \boldsymbol{\Upsilon}_0\} = 1$. Furthermore

$$\text{tr}(\mathbf{Y}' \mathbf{A} \mathbf{Y}) = \text{tr}(\mathbf{Y}'_1 \mathbf{A} \mathbf{Y}_1) + \text{tr}(\mathbf{Y}'_0 \mathbf{A} \mathbf{Y}_0). \quad (10.50)$$

Corollary 10.14.2 The distribution of the trace and other properties of $\mathbf{S} \sim \mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \boldsymbol{\Delta})$ can be described by choosing $\mathbf{A} = \mathbf{I}_N$. In particular, if $\mathbf{Y} = (\mathbf{S})\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \boldsymbol{\Sigma})$, with $\boldsymbol{\Sigma} = \boldsymbol{\Upsilon} \text{Dg}(\boldsymbol{\lambda}) \boldsymbol{\Upsilon}'$ and $\boldsymbol{\Upsilon} = [\boldsymbol{\Upsilon}_1 \ \boldsymbol{\Upsilon}_0]$, then

$$\text{tr}(\mathbf{Y}' \mathbf{Y}) = \sum_{k=1}^{p_1} \lambda_k x_k + \text{tr}(\boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{M} \boldsymbol{\Upsilon}_0), \quad (10.51)$$

with independent $x_k \sim \chi^2(\nu, \mathbf{m}'_k \mathbf{m}_k / \lambda_k)$. If $p_1 = p$ then $\text{tr}(\boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{M} \boldsymbol{\Upsilon}_0) = 0$. If $\mathbf{M} = \mathbf{0}$, then $\text{tr}(\boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{M} \boldsymbol{\Upsilon}_0) = 0$ and $x_k \sim \chi^2(\nu, 0)$. If $p_1 = p$ or $\mathbf{M} = \mathbf{0}$, then $\text{tr}(\boldsymbol{\Upsilon}'_0 \mathbf{M}' \mathbf{M} \boldsymbol{\Upsilon}_0) = 0$ and $x_k \sim \chi^2(\nu, 0)$.

Theorem 10.15 If $\mathbf{S} \sim \mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \mathbf{0})$ and $\nu \geq p$, then $|\mathbf{S}| |\boldsymbol{\Sigma}|^{-1} \sim \prod_{j=1}^p x_j$ with $x_j \sim \chi^2(\nu + 1 - j) \perp\!\!\!\perp x_{j'} \sim \chi^2(\nu + 1 - j') \forall j \neq j'$.

Proof. Here $|\mathbf{S}| |\boldsymbol{\Sigma}|^{-1} = |\mathbf{S}| |\boldsymbol{\Phi} \boldsymbol{\Phi}'|^{-1} = |\mathbf{S}| (|\boldsymbol{\Phi}| |\boldsymbol{\Phi}'|)^{-1} = |\mathbf{S}| |\boldsymbol{\Phi}|^{-1} |\boldsymbol{\Phi}'|^{-1} = |\mathbf{S}| |\boldsymbol{\Phi}^{-1}| |\boldsymbol{\Phi}^{-t}| = |\boldsymbol{\Phi}^{-t}| |\mathbf{S}| |\boldsymbol{\Phi}^{-1}| = |\boldsymbol{\Phi}^{-1} \mathbf{S} \boldsymbol{\Phi}^{-t}|$. By a previous theorem, $|\boldsymbol{\Phi}^{-1} \mathbf{S} \boldsymbol{\Phi}^{-t}| \sim \mathcal{W}_p(\nu, \mathbf{I}, \mathbf{0})$. The Bartlett decomposition theorem and corollary give $|\boldsymbol{\Phi}^{-1} \mathbf{S} \boldsymbol{\Phi}^{-t}| = |\mathbf{T} \mathbf{T}'| = |\mathbf{T}| |\mathbf{T}'| = |\mathbf{T}|^2$. Since \mathbf{T} is triangular, its determinant is the product of its diagonals,

$$|\mathbf{T}|^2 = \left(\prod_{j=1}^p t_{jj} \right)^2 = \prod_{j=1}^p t_{jj}^2. \quad (10.52)$$

The proof is completed by describing the joint distribution of $\{t_{jj}^2\}$ in terms of the distribution of $\{t_{jj}\}$ given in the triangular decomposition. \square

EXERCISES

10.1 Prove that if $\mathbf{S} = \begin{bmatrix} \mathbf{S}_{11} & \mathbf{S}_{12} \\ \mathbf{S}_{21} & \mathbf{S}_{22} \end{bmatrix} \sim \mathcal{W}_p(\nu, \boldsymbol{\Sigma}, \boldsymbol{\Delta})$ with $\mathbf{S}_{11}(q \times q)$ and $\boldsymbol{\Sigma}$ and $\boldsymbol{\Delta}$ are partitioned to match, then $\mathbf{S}_{11} \sim \mathcal{W}_q(\nu, \boldsymbol{\Sigma}_{11}, \boldsymbol{\Delta}_{11})$.

10.2 Prove directly, without using the theorem for which the following is a (very) special case:

If $\mathbf{Y} \sim \mathcal{N}_{n,p}(\mathbf{0}, \mathbf{I}_n \gamma, \mathbf{I}_p)$, then $\text{tr}(\mathbf{Y}'\mathbf{Y})/\gamma \sim \chi^2(np)$.

10.3 Consider the following SAS/IML program.

```
START GAUSS1 (N, MUMAT, FSIGMAT, SEED) ;
*Function returns matrix Gaussian N (MUMAT, I (N), SIGMA) ;
*SIGMA=FSIGMAT`*FSIGMAT;
*SEED is the random # generator initial value;
Y=MUMAT + NORMAL (J (N, NCOL (FSIGMAT), SEED)) *FSIGMAT;
RETURN (Y) ;
FINISH GAUSS1;
```

The $J(r, c, x)$ function in IML creates a matrix with r rows, c columns, and all elements equal to x . Also, $\mathbf{A}'\mathbf{B}$ indicates matrix \mathbf{A} transposed and multiplied by \mathbf{B} .

10.3.1 If \mathbf{SIGMA} is $p \times p$, what must the dimension of \mathbf{MUMAT} be?

10.3.2 What will be the dimension of the matrix \mathbf{Y} ?

10.3.3 Fully specify the distribution of data created by $\text{NORMAL}(J(N, \text{NCOL}(\text{FSIGMAT}), \text{SEED}))$. The Normal function returns pseudo-random $\mathcal{N}(0, 1)$ values. In SAS/IML, as with many pseudo-random number generators, the seed value is essentially ignored except for the very first call. Furthermore successive invocations create independent values.

10.3.4 Indicate which theorem is being used to compute \mathbf{Y} and use the theorem to fully specify the distribution of \mathbf{Y} . The program may be used to help answer the following exercise.

10.4 There are two common ways to generate a Wishart matrix in simulations. The simplest way is to generate Gaussian data and compute the Wishart as a function of the Gaussian data. The second approach is often much faster and uses the Bartlett decomposition. Here the simpler method is acceptable.

10.4.1 Generate 100 pseudo-random samples following the $\mathcal{W}_3(r, \mathbf{\Sigma}, \mathbf{0})$ distribution with $r = 10$, and $\mathbf{\Sigma} = \sigma^2[(1 - \rho)\mathbf{I}_3 + \rho\mathbf{1}_3\mathbf{1}'_3]$, with $\mathbf{1}_m$ an $m \times m$ matrix filled with 1's.

Assume $\sigma^2 = 2.0$ and $\rho = 0.50$. The result should be a set of random matrices that follow the specified Wishart distribution: $\{\mathbf{S}_i : \mathbf{S}_i \sim \mathcal{W}_3(r, \mathbf{\Sigma}, \mathbf{0})\}$ for $i \in \{1, \dots, 100\}$.

10.4.2 Print $\mathbf{\Sigma}$, \mathbf{S}_1 , \mathbf{S}_2 , \mathbf{S}_3 , and a copy of your program.

10.4.3 Compute the largest eigenvalue λ_{1i} of each of the $N = 100$ random matrices.

10.4.4 Report the mean and variance of the set of largest eigenvalues.

10.4.5 Display a frequency histogram for the distribution of the set of largest eigenvalues.

CHAPTER 11

Estimation for Univariate and Weighted Linear Models

11.1 MOTIVATION

Linear model estimation theory ranks as among the most beautiful in statistics. The derivations illustrate classical techniques for obtaining good estimators. The resulting explicit linear functions of the data also have nice geometric interpretations, which helps understand the estimators. We focus on the univariate linear model in the present chapter and leave multivariate generalizations to the next chapter. Separate treatments allows seeing the proof techniques in two slightly different settings. Although the approach creates some redundancy, taking more but shorter steps gives a faster path to understanding the multivariate GLM.

The present chapter centers on deriving estimators for β and σ^2 with a variety of optimal properties. Most of the properties do not require any particular distribution for the responses and are exact, even in small samples. In contrast, testing hypotheses uses test statistics. Finding exact test statistic distributions for small samples usually requires explicitly specifying the data distribution.

11.2 STATEMENT OF THE PROBLEM

Definition 11.1 (a) The vector β ($q \times 1$) and scalar σ^2 are the *primary parameters* of a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$, with or without Gaussian errors.

(b) Any (finite) known constant \mathbf{C} ($a \times q$) and (finite) known constant θ_0 ($a \times 1$) define *secondary parameter* $\theta = \mathbf{C}\beta + \theta_0$ ($a \times 1$).

Definition 11.2 (a) *Estimators* of primary parameters with good properties, especially unbiasedness, are indicated by $\hat{\beta}$ and $\hat{\sigma}^2$, while ones with distinct and possibly fewer desirable properties are indicated by $\tilde{\beta}$ and $\tilde{\sigma}^2$.

(b) Estimators of secondary parameters take the form $\hat{\theta} = \mathbf{C}\hat{\beta} + \theta_0$, or $\tilde{\theta} = \mathbf{C}\tilde{\beta} + \theta_0$, or $\hat{\theta} = \mathbf{C}\tilde{\beta} + \theta_0$. The third form is used only when $\hat{\theta}$ possesses a desirable property not shared by $\tilde{\beta}$, such as unbiasedness.

(c) Covariance matrices such as $\mathcal{V}(\hat{\theta})$ are also secondary parameters.

Overall, four classes of $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ are defined by allowing either FR or LTFR designs, combined with either Gaussian or unspecified distributions for responses. For $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$, with or without Gaussian errors, estimation theory for primary parameter $\boldsymbol{\beta}$ and secondary parameter $\boldsymbol{\theta}$ is closely tied to $r = \text{rank}([\mathbf{X}' \mathbf{R}'']')$.

The distinction between FR and LTFR describes properties of the columns of $[\mathbf{X}' \mathbf{R}'']'$, which correspond directly to properties of the rows of $\boldsymbol{\beta}$. In the absence of restrictions on $\boldsymbol{\beta}$ the classification depends solely on $r = \text{rank}(\mathbf{X})$. Full rank models have an unbiased estimator for $\boldsymbol{\beta}$ and for all $\boldsymbol{\theta}$, while LTFR models never have an unbiased estimator for $\boldsymbol{\beta}$ and have an unbiased estimator only for some $\boldsymbol{\theta}$. Whether or not the model is FR, unbiased estimators of σ^2 and related properties are always available [as long as $N > r = \text{rank}([\mathbf{X}' \mathbf{R}'']')$].

11.3 (UNRESTRICTED) LINEARLY EQUIVALENT LINEAR MODELS

The basic theory of transformations between linearly equivalent univariate linear models is contained in the present section and in Section 11.15. The results allow defining two important types of linear equivalence, namely (1) equivalence between a LTFR and FR model and (2) equivalence between an explicitly and an implicitly restricted model. The value of the concept lies in being able to work with a model that is simpler (often due to involving fewer parameters and variables) while being assured of not losing access to any information from the original model.

Definition 11.3 $GLM_{N,q_1}(y_i; \mathbf{X}_{i,1}\boldsymbol{\beta}_1, \sigma^2)$ and $GLM_{N,q_2}(y_i; \mathbf{X}_{i,2}\boldsymbol{\beta}_2, \sigma^2)$ are *linearly equivalent* whenever (1) for any $\boldsymbol{\beta}_1$ there exists $\boldsymbol{\beta}_2$ such that $\mathbf{X}_1\boldsymbol{\beta}_1 = \mathbf{X}_2\boldsymbol{\beta}_2$ and (2) for any $\boldsymbol{\beta}_2$ there exists $\boldsymbol{\beta}_1$ such that $\mathbf{X}_1\boldsymbol{\beta}_1 = \mathbf{X}_2\boldsymbol{\beta}_2$.

Linear equivalence describes the expected values, the means, of $\{y_i\}$, because $E(\mathbf{y}) = \mathbf{X}_1\boldsymbol{\beta}_1$ and $E(\mathbf{y}) = \mathbf{X}_2\boldsymbol{\beta}_2$. The definition implicitly requires \mathbf{X}_1 and \mathbf{X}_2 have the same number of rows and $\text{rank}(\mathbf{X}_1) = \text{rank}(\mathbf{X}_2)$. For a given $\boldsymbol{\beta}_1$ the required $\boldsymbol{\beta}_2$ need not be unique, while, given $\boldsymbol{\beta}_2$, the required $\boldsymbol{\beta}_1$ need not be unique.

Lemma 11.1 (a) Models $GLM_{N,q_1}(y_i; \mathbf{X}_{i,1}\boldsymbol{\beta}_1, \sigma^2)$ and $GLM_{N,q_2}(y_i; \mathbf{X}_{i,2}\boldsymbol{\beta}_2, \sigma^2)$ are linearly equivalent if and only if \mathbf{X}_1 and \mathbf{X}_2 (1) have N rows and (2) their columns span the same subspace of \mathfrak{R}^N .

(b) Alternately, $GLM_{N,q_1}(y_i; \mathbf{X}_{i,1}\boldsymbol{\beta}_1, \sigma^2)$ and $GLM_{N,q_2}(y_i; \mathbf{X}_{i,2}\boldsymbol{\beta}_2, \sigma^2)$ are linearly equivalent if and only if (1) $\mathbf{X}_2 = \mathbf{X}_1\mathbf{T}_1$ and (2) $\mathbf{X}_1 = \mathbf{X}_2\mathbf{T}_2$.

Proof. Left as an exercise.

A clearer understanding of equivalence may be achieved by considering the SVDs of the two design matrices. Necessarily the number of rows and ranks of \mathbf{X}_1 and \mathbf{X}_2 must be the same (otherwise one would contain information not present in the other). With $\text{Dg}(\mathbf{s}_{1,j})$ containing only the strictly positive singular values of dimension $\text{rank}(\mathbf{X}_j)$, it helps to write the SVDs (using Lemma 1.29) as

$$\begin{aligned} \mathbf{X}_j &= \mathbf{L}_j \begin{bmatrix} \text{Dg}(\mathbf{s}_j) \\ \mathbf{0} \end{bmatrix} \mathbf{R}'_j \\ &= \mathbf{L}_{1,j} \text{Dg}(\mathbf{s}_{1,j}) \mathbf{R}'_{1,j} \end{aligned} \quad (11.1)$$

and define

$$\begin{aligned} \beta_{1,j} &= \mathbf{L}'_{1,j} \mathbf{X}_j \beta_j \\ &= \text{Dg}(\mathbf{s}_{1,j}) \mathbf{R}'_{1,j} \beta_j. \end{aligned} \quad (11.2)$$

Here the 1 subscript indicates the component of the original matrix corresponding to positive singular values and the full-rank basis. Although β_1 and β_2 may have different dimensions, $\beta_{1,1}$ and $\beta_{1,2}$ must have the same dimension, because otherwise \mathbf{X}_1 and \mathbf{X}_2 would have different ranks. Requiring that for any β_1 there exists β_2 such that $\mathbf{X}_1 \beta_1 = \mathbf{X}_2 \beta_2$ implies

$$\mathbf{X}_1 \beta_1 = \mathbf{X}_2 \beta_2 \quad (11.3)$$

$$\mathbf{L}_{1,1} \text{Dg}(\mathbf{s}_{1,1}) \mathbf{R}'_{1,1} \beta_1 = \mathbf{L}_{1,2} \text{Dg}(\mathbf{s}_{1,2}) \mathbf{R}'_{1,2} \beta_2 \quad (11.4)$$

$$\mathbf{L}_{1,1} \beta_{1,1} = \mathbf{L}_{1,2} \beta_{1,2}. \quad (11.5)$$

If \mathbf{X}_1 and \mathbf{X}_2 span the same column space, then $\mathbf{L}_{1,1} = \mathbf{L}_{1,2} \mathbf{T}$ with \mathbf{T} square and FR. Therefore

$$\mathbf{L}_{1,2} \mathbf{T} \beta_{1,1} = \mathbf{L}_{1,2} \beta_{1,2} \quad (11.6)$$

$$\mathbf{L}'_{1,2} \mathbf{L}_{1,2} \mathbf{T} \beta_{1,1} = \mathbf{L}'_{1,2} \mathbf{L}_{1,2} \beta_{1,2} \quad (11.7)$$

$$\mathbf{T} \beta_{1,1} = \beta_{1,2}. \quad (11.8)$$

In turn, any parameters corresponding to a basis of one design matrix are full rank transformations of corresponding parameters in a linearly equivalent model (even though β_1 and β_2 need not be).

The definition does not directly address the possibility of equivalence between an unrestricted model (1) and a restricted model (2) or equivalence between two restricted models. The first case is covered by results in Section 11.13 which allows finding an unrestricted model (3) which is linearly equivalent to model 2. The question is then reduced to comparing models 1 and 3. Similarly, two restricted models may be compared by first transforming each separately to an unrestricted model and then comparing the results.

The guaranteed existence of equivalencies between FR and LTFR models and between restricted and unrestricted models allows deriving many of the desired results in the unrestricted and FR case. Naturally care must be taken in differentiating between properties that do generalize (from FR to LTFR and from unrestricted to restricted) from properties that do not.

Traditionally, the most important transformations between *linearly equivalent* linear models were from LTFR models to FR models. The transformations allow computing in the FR setting rather than in the LTFR setting. The following theorem formally establishes that one is never required to work with a LTFR model. One can always transform to a linearly equivalent FR model and do all estimation (theory and data analysis) in the FR setting.

Theorem 11.1 Every LTFR model has a linearly equivalent FR model. In particular, for $\text{GLM}_{N,q}\text{LTFR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ and $r = \text{rank}(\mathbf{X})$ there exists $\text{GLM}_{N,r}\text{FR}(y_i; \mathbf{X}_{i1}\boldsymbol{\beta}_1, \sigma^2)$ which is linearly equivalent. The FR model provides a *reparameterization* of the LTFR.

Proof. With the subscript 1 indicating the components corresponding to positive singular values, the SVD results summarized in Lemma 1.29 give

$$\begin{aligned} \mathbf{X} &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{R}'_1 \\ \mathbf{R}'_0 \end{bmatrix} \\ &= \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \mathbf{R}'_1, \end{aligned} \tag{11.9}$$

which allows defining

$$\mathbf{X}_1 = \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \tag{11.10}$$

$$\boldsymbol{\beta}_1 = \mathbf{R}'_1 \boldsymbol{\beta}. \tag{11.11}$$

Furthermore

$$\begin{aligned} \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\ &= \mathbf{X}_1\boldsymbol{\beta}_1 + \mathbf{e}, \end{aligned} \tag{11.12}$$

with \mathbf{X}_1 having full rank. Although especially elegant and well behaved, the model is merely one choice among infinitely many linearly equivalent models. \square

11.4 ESTIMABILITY AND CRITERIA FOR CHECKING IT

Definition 11.4 In a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ or $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Upsilon})$, $\boldsymbol{\beta}$ is the primary expected-value parameter while $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is a secondary (expected-value) parameter, as is $\boldsymbol{\theta} + \boldsymbol{\theta}_0$.

(a) Primary parameter $\boldsymbol{\beta}$ is *estimable* if and only if a $q \times N$ constant matrix \mathbf{A}_1 exists such that $\text{E}(\mathbf{A}_1\mathbf{y}) = \boldsymbol{\beta}$.

(b) Secondary parameter $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is *estimable* if and only if a constant matrix \mathbf{A}_2 ($a \times N$) exists such that $\text{E}(\mathbf{A}_2\mathbf{y}) = \boldsymbol{\theta}$.

(c) For known, fixed $\boldsymbol{\theta}_0$, $\boldsymbol{\theta} + \boldsymbol{\theta}_0$ is *estimable* if and only if $\boldsymbol{\theta}$ is *estimable*.

Part (a) is essentially redundant since it is a special case of part (b), with $\boldsymbol{\beta}$ a secondary parameter and $\mathbf{C} = \mathbf{I}_q$.

Although estimability and parameter definition stem from different issues, Helms (1988a) proved a secondary parameter is estimable if and only if it is well defined. In the present chapter, understanding *estimability*, apart from the issue of being well defined, suffices for the study of estimation. The next theorem leads to the conclusion that nonestimability is a problem *only* in LTFR models.

Theorem 11.2 Primary parameter β and secondary parameter $\theta = C\beta$, defined by known constant C may or may not be estimable in $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$.

- (a) If $\text{rank}(X) = q$, then β and θ are always estimable in $GLM_{N,q}FR(y_i; X_i\beta, \sigma^2)$.
- (b) If $\text{rank}(X) = r < q$, then β is never estimable in $GLM_{N,q}LTFR(y_i; X_i\beta, \sigma^2)$, while θ may or may not be estimable.

Proof of (a). The FR assumption $q = \text{rank}(X) = \text{rank}(X'X)$ implies $(X'X)^{-1}$ exists. If $A_1 = (X'X)^{-1}X'$ ($q \times N$) and $A_2 = CA_1$ ($a \times N$), then $\hat{\beta} = A_1y$ and $\hat{\theta} = A_2y$ are linear estimators, which gives $E(\hat{\beta}) = A_1E(y) = (X'X)^{-1}X'X\beta = \beta$ and $E(\hat{\theta}) = CA_1E(y) = C\beta = \theta$. We have described linear unbiased estimators. Therefore both β and θ are *estimable*.

Proof of (b). By contradiction. If A exists defining estimator $\hat{\beta} = Ay$ such that $E(\hat{\beta}) = E(Ay) = AX\beta = \beta \quad \forall \beta$, then $AX = I_q$, which implies $\text{rank}(AX) = q$. However, $\text{rank}(AX) \leq \min\{\text{rank}(A), \text{rank}(X)\} \leq r < q$. The assumption A exists has led to a contradiction. Therefore no such A exists. \square

Theorem 11.3 If constant C defines $\theta = C\beta$ for $GLM_{N,q}LTFR(y_i; X_i\beta, \sigma^2)$, then θ is estimable, i.e., T exists such that $E(Ty) = \theta$,

- $\Leftrightarrow A$ exists such that $C = AX$,
- \Leftrightarrow rows of C are linear combinations of rows of X , and
- \Leftrightarrow rows of C are in the space spanned by the rows of X .

Proof. (\Rightarrow) If A exists such that $C = AX$ then the estimator defined by $\hat{\theta} = Ay$ is unbiased, $E(\hat{\theta}) = E(Ay) = AX\beta = C\beta = \theta$.

(\Leftarrow) If θ is estimable, then, by definition, matrix T exists such that $E(TY) = TX\beta = \theta = C\beta$ for any β . Thus $TX\beta = C\beta \quad \forall \beta$ implies $TX = C$ and $A = T$ exists, as required, although is not necessarily unique. \square

For checking estimability, the preceding theorem has the disadvantage of involving X , which has N rows. In practice, slow and inaccurate computations can make determining estimability difficult. The following theorems use $a \times q$ matrices, which are usually much smaller, with $q \ll N$. Computationally, if $(X'X)^{-}$ is available (as it typically will be), then the next theorem provides the easiest check for estimability. Computer arithmetic reduces the task to checking $\max\{\text{abs}[C - C(X'X)^-(X'X)]\} \leq \{\max[\text{abs}(C)]\}\epsilon$ for ϵ a value judged to be numeric zero, such as 10^{-12} in many contemporary software environments. The $\text{abs}()$ operator gives the matrix of absolute values of the original elements. The

next theorem also simplifies derivations of moments and distributions of estimators in the univariate and multivariate GLM. The subsequent two theorems can serve the same purposes, although perhaps less conveniently.

Theorem 11.4 If $\text{GLM}_{N,q}\text{LTFR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ has secondary parameter $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$, with $(\mathbf{X}'\mathbf{X})^-$ any particular generalized inverse, then $\boldsymbol{\theta}$ is estimable $\Leftrightarrow \mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^-(\mathbf{X}'\mathbf{X})$.

Proof. (\Rightarrow) If $\boldsymbol{\theta}$ is estimable, then constant \mathbf{A} ($a \times N$) exists such that $\mathbf{C} = \mathbf{A}\mathbf{X}$. In turn

$$\begin{aligned} \mathbf{C} [(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X}] &= \mathbf{A}\mathbf{X}[(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X}] \\ &= \mathbf{A}[\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X}] \\ &= \mathbf{A}\mathbf{X} = \mathbf{C}. \end{aligned} \tag{11.13}$$

Searle (1971, p. 20) proved $\mathbf{X} = \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X}$, which is in Theorem 1.15.

(\Leftarrow) Given $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^- (\mathbf{X}'\mathbf{X})$, if $\mathbf{A} = \mathbf{C}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$, then $\mathbf{C} = \mathbf{A}\mathbf{X}$ and $\boldsymbol{\theta}$ is estimable by Theorem 11.3. \square

Corollary 11.4 Secondary parameter $\boldsymbol{\theta}$ is estimable $\Leftrightarrow \mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^+(\mathbf{X}'\mathbf{X})$.

Proof. (\Rightarrow) Estimable implies $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^- (\mathbf{X}'\mathbf{X})$, which clearly implies $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^+(\mathbf{X}'\mathbf{X})$.

(\Leftarrow) If $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^+(\mathbf{X}'\mathbf{X})$, then $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^- (\mathbf{X}'\mathbf{X})$. The result follows by the theorem. \square

Theorem 11.5 A $\text{GLM}_{N,q}\text{LTFR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ may have $N \times (r + s)$ \mathbf{X} partitioned as $\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2]$ with $N \times r$ \mathbf{X}_1 and $N \times s$ \mathbf{X}_2 , while $s = q - r$. If $(\mathbf{X}'_1\mathbf{X}_1)$ is nonsingular and $\mathbf{C} = [\mathbf{C}_1 \mathbf{C}_2]$ with \mathbf{C}_1 $a \times r$ and \mathbf{C}_2 $a \times s$, then $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is estimable if and only if $\mathbf{C}_2 = \mathbf{C}_1(\mathbf{X}'_1\mathbf{X}_1)^{-1}\mathbf{X}'_1\mathbf{X}_2$.

Proof. (Roy, 1957). Left as an exercise.

Since $\text{rank}(\mathbf{X}) = r$, \mathbf{X} has r linearly independent columns. It may be necessary to permute the columns of \mathbf{X} to make the columns the first r columns. The rows of $\boldsymbol{\beta}$ and then \mathbf{C} must be permuted to match.

Theorem 11.6 If $\text{GLM}_{N,q}\text{LTFR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ has $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$, then $\boldsymbol{\theta}$ is estimable $\Leftrightarrow \exists \mathbf{D}$ ($a \times q$) such that $\mathbf{C} = \mathbf{D}(\mathbf{X}'\mathbf{X})$. If \mathbf{D} exists it is unique.

Proof. (\Rightarrow) If $\mathbf{D}_1 = \mathbf{C}(\mathbf{X}'\mathbf{X})^-$ ($a \times q$), with $(\mathbf{X}'\mathbf{X})^-$ any particular generalized inverse, and $\boldsymbol{\theta}$ is estimable, then $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^- (\mathbf{X}'\mathbf{X})$, which allows writing $\mathbf{C} = \mathbf{D}_1(\mathbf{X}'\mathbf{X})$ and ensures \mathbf{D} exists.

(\Leftarrow) If \mathbf{D}_2 exists such that $\mathbf{C} = \mathbf{D}_2(\mathbf{X}'\mathbf{X})$, then $\mathbf{A} = \mathbf{D}_2\mathbf{X}'$ exists such that $\mathbf{C} = \mathbf{A}\mathbf{X}$. Therefore $\boldsymbol{\theta}$ is estimable.

Proving uniqueness is left as an exercise. \square

Example 11.1 Two different model formulations for the same simple situation illustrate some of the complications arising in considering estimability for LTFR models. For $N = 3$, two experimental units receive treatment 1, and one experimental unit receives treatment 2. Model 1, a cell-mean model, with $q = 2$, has scalar equation $y_{ij} = \mu_i + e_{ij}$ and corresponds to a $GLM_{N,2}FR(y_i; \mathbf{X}_{i,1}\beta_1, \sigma^2)$ with Gaussian errors. Model 2, a classical ANOVA model, with $q = 3$, has scalar equation $y_{ij} = \mu + \alpha_i + e_{ij}$ and corresponds to a $GLM_{N,3}LTFR(y_i; \mathbf{X}_{i,2}\beta, \sigma^2)$ with Gaussian errors. In the setting described,

$$\begin{aligned}
 y_{ig} &= \text{observation } i \text{ subject to treatment } g \\
 \mu_g &= E(y_{ig}) \\
 \mu &= \text{an "overall mean"} \\
 \alpha_g &= \text{"the effect of treatment } g\text{"} \\
 e_{ig} &\sim \mathcal{N}_1(0, \sigma^2) \text{ i.i.d.}
 \end{aligned}$$

Also

$$\begin{aligned}
 \mathbf{y} &= \mathbf{X}_1\beta_1 + \mathbf{e} \\
 \begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \end{bmatrix} &= \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} e_{11} \\ e_{12} \\ e_{21} \end{bmatrix}
 \end{aligned} \tag{11.14}$$

and

$$\begin{aligned}
 \mathbf{y} &= \mathbf{X}_2\beta_2 + \mathbf{e} \\
 \begin{bmatrix} y_{11} \\ y_{12} \\ y_{21} \end{bmatrix} &= \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{bmatrix} + \begin{bmatrix} e_{11} \\ e_{12} \\ e_{21} \end{bmatrix}.
 \end{aligned} \tag{11.15}$$

It is easily verified that

$$\hat{\beta}_1 = \begin{bmatrix} \bar{y}_1 \\ y_{21} \end{bmatrix} = \begin{bmatrix} \hat{\mu}_1 \\ \hat{\mu}_2 \end{bmatrix} \tag{11.16}$$

is unbiased, $E(\hat{\beta}_1) = \beta_1$. Furthermore any linear combination of μ_1 and μ_2 , say $C\beta_1 = c_1\mu_1 + c_2\mu_2$, is estimable since $C\hat{\beta}_1$ is a linear unbiased estimator of $C\beta_1$.

In contrast, β_2 is not estimable (no linear unbiased estimator of β_2 exists), and element j of β_2 is not estimable. In particular, no individual linear estimators with expectation equal to μ , or to α_1 , or to α_2 exist. They are not individually estimable, even as secondary parameters ($\theta = \mathbf{c}_j\beta_2 = \beta_{j,2}$, with $\mathbf{c}_j \mathbf{1} \times q$ with all zero elements except for a 1 at position j).

For model 2 the normal equations are

$$\begin{aligned}
 (\mathbf{X}'_2\mathbf{X}_2)\widehat{\boldsymbol{\beta}}_2 &= \mathbf{X}'_2\mathbf{y} \\
 &= \begin{bmatrix} 3 & 2 & 1 \\ 2 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \widehat{\mu} \\ \widehat{\alpha}_1 \\ \widehat{\alpha}_2 \end{bmatrix} \\
 &= \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{y}.
 \end{aligned} \tag{11.17}$$

Subtracting multiples of the third row from the first two rows gives the equivalent equations

$$\begin{bmatrix} 0 & 2 & -2 \\ 0 & 2 & -2 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \widehat{\mu} \\ \widehat{\alpha}_1 \\ \widehat{\alpha}_2 \end{bmatrix} = \begin{bmatrix} 1 & 1 & -2 \\ 1 & 1 & -2 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{y}. \tag{11.18}$$

Since two of the equations are identical, we have a system of two equations in three unknowns. Infinitely many solutions exist for such a system. If $\widehat{\boldsymbol{\beta}}_2$ is a solution, then so is $\widetilde{\boldsymbol{\beta}}_2 = \widehat{\boldsymbol{\beta}}_2 + \mathbf{z}$ for any $\mathbf{z} = [-z_0 \ z_0 \ z_0]'$, because

$$\mathbf{X}_2\mathbf{z} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} -z_0 \\ z_0 \\ z_0 \end{bmatrix} = \mathbf{0}. \tag{11.19}$$

Particular solutions are obtained by placing an additional linear restriction on $\boldsymbol{\beta}_2$, specifically $\mathbf{r}'\boldsymbol{\beta}_2 = a$, to supply a third equation. Requiring $\mu = 0$ corresponds to $\mathbf{r}' = [1 \ 0 \ 0]$ and $a = 0$. The choice implies three equations in three unknowns,

$$\begin{bmatrix} 1 & 0 & 0 \\ 0 & 2 & -2 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \widehat{\mu} \\ \widehat{\alpha}_1 \\ \widehat{\alpha}_2 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & -2 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{y}. \tag{11.20}$$

Further manipulation yields the solution

$$\widehat{\boldsymbol{\beta}}_2 = \begin{bmatrix} \widehat{\mu} \\ \widehat{\alpha}_1 \\ \widehat{\alpha}_2 \end{bmatrix} = \begin{bmatrix} 0 \\ (y_{11} + y_{12})/2 \\ y_{21} \end{bmatrix}. \tag{11.21}$$

It is interesting to look at the expected values of the estimators. They are

$$E\left(\widehat{\boldsymbol{\beta}}_2\right) = E\left(\begin{bmatrix} \widehat{\mu} \\ \widehat{\alpha}_1 \\ \widehat{\alpha}_2 \end{bmatrix}\right) = \begin{bmatrix} 0 \\ \mu + \alpha_1 \\ \mu + \alpha_2 \end{bmatrix} \neq \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{bmatrix} = \boldsymbol{\beta}_2. \tag{11.22}$$

We can also examine the expected value of all other solutions since they must be of the form $\widetilde{\boldsymbol{\beta}}_2 = \widehat{\boldsymbol{\beta}}_2 + \mathbf{z}$. No known quantity z_0 exists such that

$$E(\widehat{\beta}_2 + \mathbf{z}) = \begin{bmatrix} 0 \\ \mu + \alpha_1 \\ \mu + \alpha_2 \end{bmatrix} + E\left(\begin{bmatrix} -z_0 \\ z_0 \\ z_0 \end{bmatrix}\right) = \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{bmatrix} = \beta_2. \quad (11.23)$$

Regardless of model specifications, $\mu_{ij} = E(y_{ij})$ is always estimable, which is a trivial result since y_{ij} is a linear unbiased estimator of $\mu_{ij} = E(y_{ij})$. In the example, the fundamental parameters are $E(y_{11}) = \mu + \alpha_1 = E(y_{12})$ and $E(y_{21}) = \mu + \alpha_2$. Since the parameters are estimable, linear combinations of them are also estimable, including $\alpha_1 - \alpha_2 = (\mu + \alpha_1) - (\mu + \alpha_2)$. However, no linear combination of the fundamental estimable parameter has expectation equal to μ , or to α_1 , or to α_2 . They are not estimable, even as secondary parameters. The proofs of the assertions are recommended as exercises.

In spite of the estimability problems, LTFR models are popular and can be useful in some circumstances. If one restricts interest to estimable secondary parameters, such parameters may be estimated equally well via either FR or LTFR models (aside from certain computational problems for LTFR models). The FR models have an advantage when defining a secondary parameter because one need not be concerned with whether the parameter is estimable.

Since it is embarrassing to present an estimate of a nonestimable parameter, when using LTFR models one must check each secondary parameter for estimability, which can be a nuisance. Verifying the condition $C = C(X'X)^-(X'X)$ seems the simplest.

The following theorem formally justifies taking advantage of the convenience of working with a linearly equivalent model. Analyzing a full-rank model which is linearly equivalent to a less-than-full-rank model avoids the need to check for estimability and loses no information available in the original model.

Theorem 11.7 Any primary or secondary expected-value parameter estimable in $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$ is also estimable in a linearly equivalent model.

Proof. Considering estimable $\theta = C\beta$ first allows treating β as a special case. Direct comparison of the definitions of estimable and linearly equivalent provides the basis of the result. Details are left as an exercise.

11.5 CODING SCHEMES AND THE ESSENCE MATRIX

Even simple design matrices can be coded in a variety of ways. The choice affects the definitions and estimability of expected-value parameters β and $\theta = C\beta$. Readers with limited knowledge of coding schemes will find it profitable to consult Chapters 12–16 in Muller and Fetterman (2002). They provided extensive guidance concerning practical applications in ANOVA and regression. They also explicitly describe equivalencies among models using such coding schemes in one- and two-way designs.

In the univariate linear model $GLM_{N,g}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, with corresponding model equation for all observations $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$, the design matrix \mathbf{X} represents between-subject information (distinctions among independent sampling units, ISUs). Most often, data analysts use one of six coding schemes to code distinctions among G groups of observations: reference cell (regression), cell mean, effect, classical ANOVA, natural polynomial, and orthogonal polynomial.

Example 11.2 A $G = 3$ group ANOVA design, with N_g independent observations in group g , may be written with reference cell coding as

$$\begin{aligned} \mathbf{y} &= \mathbf{X}_1\boldsymbol{\beta} + \mathbf{e} \\ &= \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{1} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} & \mathbf{1} \end{bmatrix} \boldsymbol{\beta} + \mathbf{e} \\ &= \begin{bmatrix} \mathbf{1}_{N_1} \otimes [\mathbf{1} & \mathbf{0} & \mathbf{0}] \\ \mathbf{1}_{N_2} \otimes [\mathbf{1} & \mathbf{1} & \mathbf{0}] \\ \mathbf{1}_{N_3} \otimes [\mathbf{1} & \mathbf{0} & \mathbf{1}] \end{bmatrix} \boldsymbol{\beta} + \mathbf{e}. \end{aligned} \tag{11.24}$$

If $N_1 = N_2 = N_3$ (a balanced design), then $N_g = N/G$ and

$$\mathbf{X}_1 = \mathbf{1}_{N_g} \otimes \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{1} & \mathbf{1} & \mathbf{0} \\ \mathbf{1} & \mathbf{0} & \mathbf{1} \end{bmatrix} = \mathbf{1}_{N_g} \otimes \mathbf{X}_{E,1}. \tag{11.25}$$

Cell mean coding gives

$$\begin{aligned} \mathbf{y} &= \mathbf{X}_2\boldsymbol{\beta} + \mathbf{e} \\ &= \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} \end{bmatrix} \boldsymbol{\beta} + \mathbf{e} \\ &= \begin{bmatrix} \mathbf{1}_{N_1} \otimes [\mathbf{1} & \mathbf{0} & \mathbf{0}] \\ \mathbf{1}_{N_2} \otimes [\mathbf{0} & \mathbf{1} & \mathbf{0}] \\ \mathbf{1}_{N_3} \otimes [\mathbf{0} & \mathbf{0} & \mathbf{1}] \end{bmatrix} \boldsymbol{\beta} + \mathbf{e}. \end{aligned} \tag{11.26}$$

If the design is balanced, then

$$\mathbf{X}_2 = \mathbf{1}_{N_g} \otimes \begin{bmatrix} \mathbf{1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1} \end{bmatrix} = \mathbf{1}_{N_g} \otimes \mathbf{X}_{E,2}. \tag{11.27}$$

Obviously $\mathbf{X}_{E,2} = \mathbf{I}_3$.

Definition 11.5 The *essence matrix* (Helms, 1988a, b) simplifies any discussion of coding schemes. With N observations and q predictors, an $N \times q$ design matrix, \mathbf{X} , has $G \times q$ essence matrix $\text{Es}(\mathbf{X})$, which contains one and only one copy of each unique row of \mathbf{X} .

With G rows in $\text{Es}(\mathbf{X})$ each row $\mathbf{X}_g = \text{row}_g[\text{Es}(\mathbf{X})]$ identifies one of G distinct groups of sampling units. Observations which share the same $\text{row}_i(\mathbf{X})$ and hence the same \mathbf{X}_g are described as being in the same cell, which corresponds to the same treatment combination for an experiment. With N_g observations in group g and an unbalanced design,

$$\mathbf{X} = \begin{bmatrix} \mathbf{X}_1 \otimes \mathbf{1}_{N_1} \\ \mathbf{X}_2 \otimes \mathbf{1}_{N_2} \\ \vdots \\ \mathbf{X}_G \otimes \mathbf{1}_{N_G} \end{bmatrix}. \tag{11.28}$$

A balanced design has an equal number of observations in each cell so $N_g \equiv N/G$ and $\mathbf{X} = \text{Es}(\mathbf{X}) \otimes \mathbf{1}_{N/G}$.

For a one-way ANOVA design, cell-mean coding has $\text{Es}(\mathbf{X}) = \mathbf{I}_G$. Reference-cell, effect, and polynomial coding schemes also have $G \times G$ and full-rank $\text{Es}(\mathbf{X})$ for $N \times G$ and full rank \mathbf{X} . Classical ANOVA coding has less-than-full-rank $\text{Es}(\mathbf{X}) = [\mathbf{1}_G \ \mathbf{I}_G]$, while deleting any of the last G columns creates a reference-cell coding.

Questions of parameter definition and estimability can be answered in terms of the essence matrix. The easily proven fact that $\text{rank}(\mathbf{X}) = [\text{Es}(\mathbf{X})]$ gives some hint of value in considering the essence matrix. Furthermore, linear equivalence of two models can be assessed conveniently in terms of essence matrices.

11.6 UNRESTRICTED MAXIMUM LIKELIHOOD ESTIMATION OF β

The following two lemmas support the derivation of the likelihood estimators.

Lemma 11.2 If $f(\mathbf{x})$ is a positive valued function of \mathbf{x} on \mathbb{R}^p , then \mathbf{x}_0 is the location of a local maximum of $f(\mathbf{x})$ if and only if \mathbf{x}_0 is the location of a local maximum of $\log[f(\mathbf{x})]$.

Proof. (\Rightarrow) Proving $f(\mathbf{x}_0)$ is a local maximum is equivalent to proving $\exists \delta > 0$ such that $\|\mathbf{x} - \mathbf{x}_0\| < \delta \Rightarrow f(\mathbf{x}) \leq f(\mathbf{x}_0)$. The monotonicity of the $\log(\cdot)$ function allows writing $f(\mathbf{x}) \leq f(\mathbf{x}_0) \Rightarrow \log[f(\mathbf{x})] \leq \log[f(\mathbf{x}_0)]$. Therefore we know $\|\mathbf{x} - \mathbf{x}_0\| < \delta \Rightarrow \log[f(\mathbf{x})] \leq \log[f(\mathbf{x}_0)]$.

(\Leftarrow) Proving $\log[f(\mathbf{x}_0)]$ is a local maximum is equivalent to proving $\exists \delta > 0$ such that $\|\mathbf{x} - \mathbf{x}_0\| < \delta \Rightarrow \log[f(\mathbf{x})] \leq \log[f(\mathbf{x}_0)] \Rightarrow f(\mathbf{x}) \leq f(\mathbf{x}_0)$ due to the monotonicity of $\exp(\cdot)$. \square

Lemma 11.3 For $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$, the system of equations $(\mathbf{X}'\mathbf{X})\beta = \mathbf{X}'\mathbf{y}$ is consistent (Definition 1.34) when solving for unknown β in terms of known \mathbf{X} and \mathbf{y} .

Proof. $\mathbf{a}'(\mathbf{X}'\mathbf{X}) = \mathbf{0} \Rightarrow \mathbf{a}'(\mathbf{X}'\mathbf{X})\mathbf{a} = 0 \Rightarrow \mathbf{a}'\mathbf{X}' = \mathbf{0} \Rightarrow \mathbf{a}'\mathbf{X}'\mathbf{y} = \mathbf{0}$. \square

Any model with less-than-full-rank design matrix disallows finding a unique estimate for β . Infinitely many different estimates provide a valid solution, which could each be termed a supremum likelihood estimate. However, even when a unique estimate for β cannot be found, a unique estimate of σ^2 will still exist. We follow convention in describing the estimates for $\{\beta, \sigma^2\}$ as the maximum likelihood estimates, even with a less-than-full-rank design.

Theorem 11.8 For $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with Gaussian errors and $r = \text{rank}(\mathbf{X}) \leq q$, the joint supremum (for β if $r < q$) or maximum (if $r = q$ for β ; always for σ^2) likelihood estimators of β and σ^2 are

$$\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y} \tag{11.29}$$

$$\tilde{\sigma}^2 = (\mathbf{y} - \mathbf{X}\tilde{\beta})'(\mathbf{y} - \mathbf{X}\tilde{\beta})/N = \mathbf{y}'[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}']\mathbf{y}/N. \tag{11.30}$$

Here $\tilde{\beta}$ is any solution of $(\mathbf{X}'\mathbf{X})\tilde{\beta} = \mathbf{X}'\mathbf{y}$, with infinitely many for $r < q$, and one unique solution if $r = q$ and $(\mathbf{X}'\mathbf{X})^{-} = (\mathbf{X}'\mathbf{X})^{-1}$. The value of $\tilde{\sigma}^2$ is invariant to the choice of $(\mathbf{X}'\mathbf{X})^{-}$. It is customary to use $\hat{\sigma}^2 = \tilde{\sigma}^2 N/(N - r)$, which is unbiased. If \mathbf{X} ($N \times q$) is full rank, then $(\mathbf{X}'\mathbf{X})^{-} = (\mathbf{X}'\mathbf{X})^{-1}$, and $\tilde{\beta} = \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ is unique and unbiased.

Proof. The density function of \mathbf{y} ($N \times 1$) with $\mu = \mathbf{X}\beta$, $\Sigma = \sigma^2\mathbf{I}$ is

$$f_{\mathbf{y}}(\mathbf{y}_*) = (2\pi)^{-N/2} |\Sigma|^{-1/2} \exp[-(\mathbf{y}_* - \mu)' \Sigma^{-1} (\mathbf{y}_* - \mu)/2]. \tag{11.31}$$

We actually maximize the log likelihood (using Lemma 11.2), which is

$$\begin{aligned} \log L(\mathbf{y}_*; \beta, \sigma^2) &= \log[f_{\mathbf{y}}(\mathbf{y}_*)] \\ &= -\frac{N}{2} \log(2\pi) - \frac{N}{2} \log(\sigma^2) - \frac{1}{2} (\mathbf{y}_* - \mathbf{X}\beta)' (\mathbf{y}_* - \mathbf{X}\beta) / \sigma^2 \\ &= c - \frac{N}{2} \log(\sigma^2) - (\mathbf{y}'_*\mathbf{y}_* - 2\mathbf{y}'_*\mathbf{X}\beta + \beta'\mathbf{X}'\mathbf{X}\beta) / (2\sigma^2). \end{aligned} \tag{11.32}$$

Step 1. Critical points are found by finding zeros of partial derivatives of $\log L$ with respect to $\tau = [\beta' \ \sigma^2]'$. Here $\partial \log L / \partial \beta = (2\sigma^2)^{-1} [2\mathbf{X}'\mathbf{y} - 2(\mathbf{X}'\mathbf{X})\beta]$ is $q \times 1$. Evaluating $\mathbf{0} = \partial \log L / \partial \beta$ at $\beta = \tilde{\beta}$ implies $(\mathbf{X}'\mathbf{X})\tilde{\beta} = \mathbf{X}'\mathbf{y}$. Equivalently, $\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}_*$ since the equations are consistent by Lemma 1.17. Also,

$$\frac{\partial \log L}{\partial \sigma^2} = -\frac{N}{2\sigma^2} + \frac{1}{2(\sigma^2)^2} [(\mathbf{y}_* - \mathbf{X}\beta)' (\mathbf{y}_* - \mathbf{X}\beta)]. \tag{11.33}$$

Evaluating $\mathbf{0} = \partial \log L / \partial \sigma^2$ at $\sigma^2 = \tilde{\sigma}^2$ implies $\tilde{\sigma}^2 = (\mathbf{y} - \mathbf{X}\tilde{\beta})'(\mathbf{y} - \mathbf{X}\tilde{\beta})/N$ or $\tilde{\sigma}^2 = \mathbf{y}'(\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}')\mathbf{y}/N$. The estimator is invariant to the choice of generalized inverse by Theorem 1.15.

Step 2. Is the critical point a maximum, minimum, or saddle point? Deciding requires proving the Hessian matrix

$$\frac{\partial^2 \log L}{\partial \tau \partial \tau'} = \begin{bmatrix} \frac{\partial^2 \log L}{\partial \sigma^2 \partial \sigma^2} & \frac{\partial^2 \log L}{\partial \sigma^2 \partial \beta'} \\ \frac{\partial^2 \log L}{\partial \beta \partial \sigma^2} & \frac{\partial^2 \log L}{\partial \beta \partial \beta'} \end{bmatrix} \quad (11.34)$$

is negative definite if $(\mathbf{X}'\mathbf{X})^{-1}$ exists and negative semidefinite otherwise. Further details are left as an exercise. \square

Lemma 11.4 The solution to the equations $(\mathbf{X}'\mathbf{X})\beta = \mathbf{X}'\mathbf{y}$ gives predicted values orthogonal to the residuals, which leads to referring to $(\mathbf{X}'\mathbf{X})\beta = \mathbf{X}'\mathbf{y}$ as the “normal” equations.

Proof.

$$\begin{aligned} \hat{\mathbf{y}}'\hat{\mathbf{e}} &= [\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}]'[\mathbf{y} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}] \\ &= \mathbf{y}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y} \\ &= \mathbf{y}'\mathbf{H}(\mathbf{I} - \mathbf{H})\mathbf{y} \\ &= \mathbf{y}'(\mathbf{H} - \mathbf{H}^2)\mathbf{y} \\ &= 0. \end{aligned} \quad (11.35) \quad \square$$

The adjective “normal” in the lemma refers to a perpendicular property and applies whether or not the data are Gaussian, as does the following lemma.

Lemma 11.5 (a) Predicted values and residuals from *any* $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ have zero covariance and zero correlation.

(b) With Gaussian errors the predicted values are statistically independent of the residuals.

Proof.

$$\begin{aligned} \mathcal{V}(\hat{\mathbf{y}}, \hat{\mathbf{e}}) &= \mathcal{V}[\mathbf{H}\mathbf{y}, (\mathbf{I} - \mathbf{H})\mathbf{y}] \\ &= \mathbf{E}\{(\mathbf{H}\mathbf{y})[(\mathbf{I} - \mathbf{H})\mathbf{y}]'\} - \mathbf{E}(\mathbf{H}\mathbf{y})\mathbf{E}\{[(\mathbf{I} - \mathbf{H})\mathbf{y}]'\} \\ &= \mathbf{H}[\mathbf{E}(\mathbf{y}\mathbf{y}')](\mathbf{I} - \mathbf{H}) - \mathbf{H}\mathbf{E}(\mathbf{y})\mathbf{E}(\mathbf{y}')(\mathbf{I} - \mathbf{H}) \\ &= \mathbf{H}[\mathbf{E}(\mathbf{y}\mathbf{y}') - \mathbf{E}(\mathbf{y})\mathbf{E}(\mathbf{y}')](\mathbf{I} - \mathbf{H}) \\ &= \mathbf{H}\mathcal{V}(\mathbf{y})(\mathbf{I} - \mathbf{H}) \\ &= \mathbf{H}\sigma^2\mathbf{I}_N(\mathbf{I} - \mathbf{H}) = \mathbf{0}_{N \times N}. \end{aligned} \quad (11.36)$$

Jointly expressing predicted values and residuals as a linear transformation of a vector Gaussian implies they are jointly vector Gaussian:

$$\begin{bmatrix} \hat{\mathbf{y}} \\ \hat{\mathbf{e}} \end{bmatrix} = \begin{bmatrix} \mathbf{H} \\ \mathbf{I} - \mathbf{H} \end{bmatrix} \mathbf{y} \sim S\mathcal{N}_{2N} \left\{ \begin{bmatrix} \mathbf{X}\beta \\ \mathbf{0} \end{bmatrix}, \sigma^2 \begin{bmatrix} \mathbf{H} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} - \mathbf{H} \end{bmatrix} \right\}. \quad (11.37)$$

Zero covariance among Gaussian vectors implies statistical independence. \square

Lemma 11.6 (a) The β estimator and the residuals from any $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ have zero covariance and zero correlation.

(b) With Gaussian errors the β estimator and the residuals are statistically independent of each other.

Proof. For $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, Theorem 1.15 gives $[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{H} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, which implies $[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'](\mathbf{I}_N - \mathbf{H}) = \mathbf{0}$. Proceeding as in the proof of the last lemma gives

$$\begin{aligned} \begin{bmatrix} \tilde{\beta} \\ \tilde{e} \end{bmatrix} &= \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ (\mathbf{I}_N - \mathbf{H}) \end{bmatrix} \mathbf{y} \\ &\sim \mathcal{SN}_{q+N} \left\{ \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta \\ \mathbf{0} \end{bmatrix}, \sigma^2 \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{I}_N - \mathbf{H}) \end{bmatrix} \right\}. \quad \square \end{aligned} \tag{11.38}$$

11.7 UNRESTRICTED BLUE ESTIMATION OF β

Definition 11.6 If at least one function of a set of data provides a linear unbiased estimator of a parameter, one may be best in the sense of having minimum variance (among unbiased estimators). Such a *best linear unbiased estimator* (BLUE) is also known as a *(linear) uniformly minimum variance unbiased estimator* (UMVUE).

Theorem 11.9 For a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i\beta, \sigma^2)$, the class of linear unbiased estimators is $\mathcal{C}_{\text{LUE}} = \{\tilde{\beta} : \tilde{\beta} = \mathbf{A}\mathbf{y}, \mathbf{A} (q \times N), \mathbf{E}(\tilde{\beta}) = \beta\}$ for \mathbf{A} constant. The BLUE can be identified in two ways. (1) The BLUE of β is $\hat{\beta} = \mathbf{A}_0\mathbf{y} \Leftrightarrow \mathcal{V}(\tilde{\beta}) - \mathcal{V}(\hat{\beta}) \equiv \mathcal{V}(\mathbf{A}\mathbf{y}) - \mathcal{V}(\mathbf{A}_0\mathbf{y})$ is *positive semidefinite* for all $\tilde{\beta} \in \mathcal{C}_{\text{LUE}}$. (2) For $\tilde{\theta}_C = \mathbf{C}\tilde{\beta}$ and $\hat{\theta}_C = \mathbf{C}\hat{\beta}$ scalar $\mathcal{V}(\tilde{\theta}_C) - \mathcal{V}(\hat{\theta}_C) \equiv \mathcal{V}(\mathbf{C}\tilde{\beta}) - \mathcal{V}(\mathbf{C}\hat{\beta})$ is nonnegative for all $\mathbf{C}' \in \mathbb{R}^q$ and $\tilde{\beta} \in \mathcal{C}_{\text{LUE}}$. Condition 1 holds \Leftrightarrow condition 2 holds.

Proof. For any $\tilde{\beta}, \hat{\beta} \in \mathcal{C}_{\text{LUE}}$, the matrix $\mathcal{V}(\tilde{\beta}) - \mathcal{V}(\hat{\beta})$ is *positive semidefinite* if and only if the scalar $\mathcal{V}(\mathbf{C}\tilde{\beta}) - \mathcal{V}(\mathbf{C}\hat{\beta})$ is nonnegative $\forall \mathbf{C}' \in \mathbb{R}^q$ because

$$\mathcal{V}(\mathbf{C}\tilde{\beta}) - \mathcal{V}(\mathbf{C}\hat{\beta}) = \mathbf{C}[\mathcal{V}(\tilde{\beta}) - \mathcal{V}(\hat{\beta})]\mathbf{C}'. \tag{11.39}$$

Matrix theory guarantees the matrix (the left side of the equation) is positive definite or positive semidefinite if and only if the quadratic form (the right side of the equation) is nonnegative $\forall \mathbf{C}$. □

Theorem 11.10 For a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i\beta, \sigma^2)$ the unique BLUE of β is $\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$. The result is known as the *Gauss-Markov theorem*.

Proof. If $\mathbf{A}_0 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ ($q \times N$), then $\hat{\boldsymbol{\beta}} = \mathbf{A}_0\mathbf{y}$ is in the class of linear unbiased estimators, $\hat{\boldsymbol{\beta}} \in \mathbb{C}_{\text{LUE}} = \{\tilde{\boldsymbol{\beta}} : \tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y} \text{ for some } \mathbf{A} \text{ and } E(\tilde{\boldsymbol{\beta}}) = \boldsymbol{\beta}\}$. An arbitrary LUE $\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y}$ has $E(\tilde{\boldsymbol{\beta}}) = E(\mathbf{A}\mathbf{y}) = \mathbf{A}\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}$ for arbitrary $\boldsymbol{\beta}$. Here $\mathbf{A}\mathbf{X}\boldsymbol{\beta} = \boldsymbol{\beta}$ for all $\boldsymbol{\beta}$ if and only if $\mathbf{A}\mathbf{X} = \mathbf{I}$. The special property \mathbf{A} ($q \times N$) must have in relation to \mathbf{X} ($N \times q$) is used to establish the fact that the matrix

$$\mathcal{V}(\tilde{\boldsymbol{\beta}}) - \mathcal{V}(\hat{\boldsymbol{\beta}}) = \mathcal{V}(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) \tag{11.40}$$

is a covariance matrix and therefore must be positive definite or positive semidefinite. It then follows $\hat{\boldsymbol{\beta}}$ is best in \mathbb{C}_{LUE} . The details are as follows. Beginning with

$$\begin{aligned} \mathcal{V}(\tilde{\boldsymbol{\beta}}) &= \mathcal{V}[(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \hat{\boldsymbol{\beta}}] \\ &= \mathcal{V}(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + \mathcal{V}(\hat{\boldsymbol{\beta}}) + 2\mathcal{V}[\hat{\boldsymbol{\beta}}, (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})] \end{aligned} \tag{11.41}$$

leads to examining the ($q \times q$) matrix

$$\begin{aligned} \mathcal{V}[\hat{\boldsymbol{\beta}}, (\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}})] &= \mathcal{V}[\mathbf{A}_0\mathbf{y}, (\mathbf{A} - \mathbf{A}_0)\mathbf{y}] \\ &= \mathbf{A}_0(\sigma^2\mathbf{I})(\mathbf{A} - \mathbf{A}_0)' \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\mathbf{A}' - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}[\mathbf{X}'\mathbf{A}' - \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{I} - \mathbf{I}) = \mathbf{0}. \end{aligned} \tag{11.42}$$

In turn, $\mathcal{V}(\tilde{\boldsymbol{\beta}}) - \mathcal{V}(\hat{\boldsymbol{\beta}}) = \mathcal{V}(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) + 2 \cdot \mathbf{0}$. Since the difference matrix is a covariance matrix, the difference matrix must be positive definite or positive semidefinite. Thus $\hat{\boldsymbol{\beta}}$ is BLUE.

The vector $\hat{\boldsymbol{\beta}}$ is also the unique BLUE, which can be proven by contradiction. Suppose $\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{y}$ and $\hat{\boldsymbol{\beta}} = \mathbf{A}_0\mathbf{y}$ are both BLUE with $\mathbf{A}_0 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \neq \mathbf{A}$. If both are BLUE, $[\mathcal{V}(\tilde{\boldsymbol{\beta}}) - \mathcal{V}(\hat{\boldsymbol{\beta}})]$ and $[\mathcal{V}(\hat{\boldsymbol{\beta}}) - \mathcal{V}(\tilde{\boldsymbol{\beta}})]$ are positive definite or positive semidefinite, which gives $[\mathcal{V}(\tilde{\boldsymbol{\beta}}) - \mathcal{V}(\hat{\boldsymbol{\beta}})] = \mathbf{0}$. Here $[\mathcal{V}(\tilde{\boldsymbol{\beta}}) - \mathcal{V}(\hat{\boldsymbol{\beta}})] = [\mathcal{V}(\hat{\boldsymbol{\beta}}) - \mathcal{V}(\tilde{\boldsymbol{\beta}})]$ implies $\mathcal{V}(\tilde{\boldsymbol{\beta}} - \hat{\boldsymbol{\beta}}) = \mathbf{0}$. In turn $\mathcal{V}[(\mathbf{A} - \mathbf{A}_0)\mathbf{y}] \equiv (\mathbf{A} - \mathbf{A}_0)\sigma^2\mathbf{I}(\mathbf{A} - \mathbf{A}_0)' = \mathbf{0}$ for all σ^2 , which implies $(\mathbf{A} - \mathbf{A}_0) = \mathbf{0}$, which is a contradiction and implies the assumption $(\mathbf{A} \neq \mathbf{A}_0)$ is false. \square

11.8 UNRESTRICTED LEAST SQUARES ESTIMATION OF $\boldsymbol{\beta}$

Here we consider building a model for a vector of observations, \mathbf{y} ($N \times 1$). Also $\hat{\mathbf{y}} = \mathbf{t}(\boldsymbol{\beta})$ models \mathbf{y} , with $\mathbf{t}(\cdot)$ a (vector-valued) transformation from $\boldsymbol{\beta} \in \mathbb{R}^q$ to $\hat{\mathbf{y}} \in \mathbb{R}^N$. Here \mathbf{B} indicates the parameter space of $\boldsymbol{\beta}$, with $\boldsymbol{\beta} \in \mathbf{B} \subset \mathbb{R}^q$.

The choice of estimator may be restricted to a certain class of estimators \mathbb{C} which satisfy certain restrictions or conditions, such as \mathbb{C}_{LUE} . Goodness of fit for such a model is assessed by examining the vector of estimated errors, the residuals $(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{r}(\tilde{\boldsymbol{\beta}}) = \hat{\mathbf{e}}$ for $\tilde{\boldsymbol{\beta}}$, an estimator of $\boldsymbol{\beta}$. The value of the residuals depend

on the choice of $\tilde{\beta}$. The study of decision theory has led to many criteria for “small error,” corresponding to many vector norms.

Definition 11.7 (a) The *squared error loss function* is $SSE(\tilde{\beta}) = \hat{e}'\hat{e} = \sum_{i=1}^N \hat{e}_i^2 = \|\hat{e}\|^2$.
(b) If $\hat{y} = t(\beta)$ is a model for y , and $\hat{\beta} \in \mathbb{C}$ for \mathbb{C} a class of estimators of $\beta \in \mathbf{B}$, then $\hat{\beta}$ is a *least squares estimator* of β if and only if $SSE(\hat{\beta}) \leq SSE(\tilde{\beta})$ for all $\tilde{\beta} \in \mathbb{C}$.

Such an estimator may or may not be unique, biased, or linear.

The Gauss-Markov theorem for the full-rank GLM provides a unique linear least-squares estimator $\hat{\beta}$ a representation which is linear in parameters, $t(\beta) = X\beta$. For each observation, the parameter estimates imply a predicted value, as in $\hat{y}_i = \hat{\beta}_1 x_{i1} + \hat{\beta}_2 x_{i2} + \hat{\beta}_3 x_{i3}$, and error estimates, the residuals, $\hat{e}_i = y_i - \hat{y}_i$. In contrast, we may seek a least squares estimator in the context of an inherently nonlinear model such as $\hat{y}_i = \hat{\beta}_1 x_{i1}^{\hat{\beta}_2} + \hat{\beta}_3 x_{i2}$.

A less-than-full-rank GLM corresponds to an inconsistent system of equations without a unique solution. Each nonunique set of parameter estimates does imply predicted values and error estimates. The following theorem, and other results later in the chapter, characterize important properties of such models.

Definition 11.8 A $GLM_{N,q}LTFR(y_i; X_i\beta, \sigma^2)$ with $r = \text{rank}(X) \leq q$, $\beta \in \mathbb{R}^q$, $\tilde{\beta} \in \mathbb{C}_{LE} \equiv \{\tilde{\beta} : \tilde{\beta} = Ay \text{ for some } A (q \times N)\}$, has $\hat{e} = y - X\tilde{\beta}$ and $SSE(\tilde{\beta}) = \hat{e}'\hat{e} = (y - X\tilde{\beta})'(y - X\tilde{\beta})$. The vector $\hat{\beta} \in \mathbb{C}_{LE}$ is a (linear) least squares estimator of β if and only if $SSE(\hat{\beta}) \leq SSE(\tilde{\beta}) \forall \tilde{\beta} \in \mathbb{C}_{LE}$.

Theorem 11.11 The following hold for $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$.

(a) Any solution $\tilde{\beta}$ of the normal equations $(X'X)\tilde{\beta} = X'y$ is a least squares estimator for β . The solutions are of the form $\tilde{\beta} = (X'X)^- X'y$ with $SSE(\tilde{\beta}) = y'[I - X(X'X)^- X']y$, which is invariant to the choice of generalized inverse.

(b) If $X (N \times q)$ has rank q , then $(X'X)^- = (X'X)^{-1}$, and $\tilde{\beta} = \hat{\beta} = (X'X)^{-1} X'y$ is unique and unbiased. Uniqueness of $\hat{\beta}$ ensures $SSE(\hat{\beta}) < SSE(\tilde{\beta})$ for all $\tilde{\beta}$ not identical to $\hat{\beta}$.

Proof. We want to minimize $SSE(b) = (y - Xb)'(y - Xb) = y'y - 2b'X'y + b'X'Xb$. Critical points are found by finding zeros of partial derivatives of SSE with respect to b . The $q \times 1$ vector of derivatives is

$$\partial SSE / \partial b = -2X'y + 2(X'X)b. \tag{11.43}$$

Evaluating $\partial SSE / \partial b = 0$ at $b = \tilde{\beta}$ implies $(X'X)\tilde{\beta} = X'y$, which implies

$\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ because the equations are consistent. Is the critical point a minimum, maximum, or saddle point? The matrix of second derivatives,

$$\frac{\partial^{(2)}SSE}{\partial\mathbf{b}\partial\mathbf{b}'} = 2\mathbf{X}'\mathbf{X}, \tag{11.44}$$

is at least positive semidefinite. Having $\mathbf{X}'\mathbf{X}$ positive definite \Rightarrow a minimum point; $\mathbf{X}'\mathbf{X}$ negative definite \Rightarrow a maximum point; $\mathbf{X}'\mathbf{X} \pm$ semidefinite \Rightarrow a minimum, maximum, or saddle point; and $\mathbf{X}'\mathbf{X}$ indefinite \Rightarrow a saddle point.

If \mathbf{X} has full rank, then $\tilde{\beta} = \hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ corresponds to a minimum and is unique and unbiased.

For the LTFR case, it must be determined whether or not $\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ corresponds to a minimum. The result can be determined by evaluating SSE for other points in the neighborhood of $\tilde{\beta}$, say $\tilde{\beta} + \mathbf{h}$, for \mathbf{h} a vector of any length in any direction. Here $SSE(\tilde{\beta} + \mathbf{h}) = \|\mathbf{y} - \mathbf{X}(\tilde{\beta} + \mathbf{h})\|^2$ can be written as

$$(\tilde{\beta} + \mathbf{h}) = \|\mathbf{y} - \mathbf{X}\tilde{\beta}\|^2 + \|\mathbf{X}\mathbf{h}\|^2 + 2\mathbf{h}'[(\mathbf{X}'\mathbf{X})\tilde{\beta} - \mathbf{X}'\mathbf{y}], \tag{11.45}$$

a Taylor series expansion. Thus $SSE(\tilde{\beta} + \mathbf{h}) = SSE(\tilde{\beta}) + \mathbf{h}'(\mathbf{X}'\mathbf{X})\mathbf{h} + 0$. Now, $\mathbf{X}'\mathbf{X}$ is positive definite $\Leftrightarrow \mathbf{h}'(\mathbf{X}'\mathbf{X})\mathbf{h} > 0 \forall \mathbf{h}$. Also, $\mathbf{X}'\mathbf{X}$ is positive semidefinite $\Leftrightarrow \mathbf{h}'(\mathbf{X}'\mathbf{X})\mathbf{h} \geq 0 \forall \mathbf{h}$ and $\mathbf{h}'(\mathbf{X}'\mathbf{X})\mathbf{h} = 0$ for some \mathbf{h} . Therefore, in the LTFR case SSE has a level valley which is minimal.

Invariance of $SSE(\tilde{\beta})$ follows from uniqueness of $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ (Theorem 1.15). □

Theorem 11.12 For a $GLM_{N,q}FR(y_i; \mathbf{X}_i\beta, \sigma^2)$ with Gaussian errors,

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \tag{11.46}$$

$$\hat{\sigma}^2 = \frac{1}{N-q}\mathbf{y}'\left[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\right]\mathbf{y} \tag{11.47}$$

are **(a)** unbiased, **(b)** consistent, **(c)** efficient, **(d)** complete, **(e)** sufficient, and **(f)** UMVUE. Furthermore **(g)** $\hat{\beta}$ and $\hat{\sigma}^2$ are mutually independent, **(h)** $\hat{\beta} \sim \mathcal{N}_q[\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]$, and **(i)** $\hat{\sigma}^2(N-q)/\sigma^2 \sim \chi^2(N-q)$.

Proof. Proofs are left to the reader.

11.9 UNRESTRICTED MAXIMUM LIKELIHOOD ESTIMATION OF θ

Least squares estimation does not apply directly to estimation of secondary parameters θ except when MLE results apply to secondary parameter estimation. Even then, some additional details must be considered.

Theorem 11.13 For $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors and $r = \text{rank}(\mathbf{X}) \leq q$, $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$, constant \mathbf{C} ($q \times q$) of rank q , the joint supremum (for $\boldsymbol{\theta}$ if $\boldsymbol{\theta}$ is not estimable) or maximum (for $\boldsymbol{\theta}$ if $\boldsymbol{\theta}$ is estimable; always for σ^2) likelihood estimators of $\boldsymbol{\theta}$ and σ^2 are

$$\tilde{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}} \tag{11.48}$$

$$\tilde{\sigma}^2 = (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})/N = \mathbf{y}'[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}']\mathbf{y}/N. \tag{11.49}$$

Here $\tilde{\boldsymbol{\beta}}$ is one of infinitely many solutions of $(\mathbf{X}'\mathbf{X})\tilde{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$ for $r < q$ and the unique solution if $r = q$ and $(\mathbf{X}'\mathbf{X})^{-} = (\mathbf{X}'\mathbf{X})^{-1}$. Estimable $\boldsymbol{\theta}$ gives $\tilde{\boldsymbol{\theta}}$ invariant to the choice of $(\mathbf{X}'\mathbf{X})^{-}$, while $\tilde{\sigma}^2$ is always invariant to $(\mathbf{X}'\mathbf{X})^{-}$.

Proof. The invariance of $\tilde{\boldsymbol{\theta}}$ for estimable $\boldsymbol{\theta}$ derives from two facts. First, estimability $\Rightarrow \mathbf{C} = \mathbf{A}\mathbf{X}$ for some \mathbf{A} . Second, $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'$ is invariant to choice of generalized inverse (Theorem 1.15). Here $\tilde{\boldsymbol{\theta}} = \hat{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}} = \mathbf{A}\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$. For nonestimable $\boldsymbol{\theta}$, infinitely many MLEs exist.

It remains to be proven that MLE estimators are of the stated form, namely $\tilde{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}}$ and $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$. Since matrix \mathbf{C} is nonsingular and $q \times q$, $\mathbf{E}(\mathbf{y}) \equiv \mathbf{X}\boldsymbol{\beta} = (\mathbf{X}\mathbf{C}^{-1})(\mathbf{C}\boldsymbol{\beta}) = \mathbf{Z}\boldsymbol{\theta}$ with $(N \times q)$ $\mathbf{Z} = \mathbf{X}\mathbf{C}^{-1}$ and $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$. Thus a linearly equivalent model with $\boldsymbol{\theta}$ as the primary expected-value parameter is $\text{GLM}_{N,q}(y_i; \mathbf{Z}_i\boldsymbol{\theta}, \sigma^2)$ with Gaussian errors and $r = \text{rank}(\mathbf{Z})$. The problem is thus reduced to one already solved. \square

Corollary 11.13 For a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors and $r = \text{rank}(\mathbf{X}) \leq q$, secondary parameter $\boldsymbol{\theta}_i = \mathbf{C}_i\boldsymbol{\beta} \in \mathfrak{R}^{a_i}$ for $i \in \{1, 2, \dots, t\}$ may be defined by \mathbf{C}_i ($a_i \times q$) such that $\text{rank}(\mathbf{C}_i) = a_i \leq q = \sum_{i=1}^t a_i$ and

$$\boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \\ \vdots \\ \boldsymbol{\theta}_t \end{bmatrix} = \begin{bmatrix} \mathbf{C}_1 \\ \mathbf{C}_2 \\ \vdots \\ \mathbf{C}_t \end{bmatrix} \boldsymbol{\beta} = \mathbf{C}\boldsymbol{\beta}. \tag{11.50}$$

- (a) A joint MLE of $\boldsymbol{\theta}_i$ ($a_i \times 1$) is $\tilde{\boldsymbol{\theta}}_i = \mathbf{C}_i\tilde{\boldsymbol{\beta}}$, in which $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$.
- (b) If $\boldsymbol{\theta}_i$ is estimable, then $\tilde{\boldsymbol{\theta}}_i$ is invariant to the choice of $(\mathbf{X}'\mathbf{X})^{-}$.
- (c) $\tilde{\boldsymbol{\theta}}_i$ is invariant to the choice of $\mathbf{C}_{i'}$ $\forall i' \neq i$.

Proof. The results follow from the theorem since \mathbf{C} ($q \times q$) is nonsingular. The rows of the \mathbf{C}_i are linearly independent, as are, collectively, the rows of \mathbf{C} . \square

Although the \mathbf{C}_i of interest may not span the estimation space, one can always find additional rows to create nonsingular \mathbf{C} ($q \times q$). It is then convenient to find the joint MLE of all the secondary parameters, $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_t\}$, and define the joint MLE of $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_k\}$, $k < t$, to be the same values as the MLE of all of $\boldsymbol{\theta}$ ($q \times 1$). The result is invariant to the choice of the additional rows making \mathbf{C} invertible.

If one attempts maximum likelihood estimation of $\theta_i, i \in \{1, 2, \dots, t - 1\}$ (or fewer), then the problem is essentially indeterminate. The elements of θ_t make an appearance as nuisance parameters, which makes the maximization depend on the unknown value of θ_t .

11.10 UNRESTRICTED BLUE ESTIMATION OF θ

Definition 11.9 (a) For $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$ with $r = \text{rank}(X) \leq q$, any $a \times 1$ estimable secondary parameter $\theta = C\beta$ has an associated class $\mathbb{C}_{LUE} = \{\tilde{\theta} : \tilde{\theta} = A\mathbf{y} \text{ for some } A (a \times N) \text{ and } E(\tilde{\theta}) = \theta\}$ of linear unbiased estimators.
(b) One element of \mathbb{C}_{LUE} may have minimum variance and therefore gives a BLUE.
(c) Such an estimator is also known as the (linear) *uniformly minimum variance unbiased estimator* (UMVUE).

Theorem 11.14 (a) The BLUE of θ is $\hat{\theta} = A_0\mathbf{y}$ such that $[\mathcal{V}((\tilde{\theta})) - \mathcal{V}(\hat{\theta})] \equiv [\mathcal{V}(A\mathbf{y}) - \mathcal{V}(A_0\mathbf{y})]$ is nonnegative definite $\forall \tilde{\theta} \in \mathbb{C}_{LUE}$.

(b) Equivalently, (scalar) $[\mathcal{V}(t'\tilde{\theta}) - \mathcal{V}(t'\hat{\theta})] \geq 0 \forall t \in \mathbb{R}^a$ and $\tilde{\theta} \in \mathbb{C}_{LUE}$.

Proof. For any $\tilde{\theta} \in \mathbb{C}_{LUE}$ and any $\hat{\theta} \in \mathbb{C}_{LUE}$, the matrix $[\mathcal{V}(\tilde{\theta}) - \mathcal{V}(\hat{\theta})]$ is nonnegative definite \Leftrightarrow the scalar $[\mathcal{V}(t'\tilde{\theta}) - \mathcal{V}(t'\hat{\theta})]$ is nonnegative $\forall t \in \mathbb{R}^a$ because $[\mathcal{V}(t'\tilde{\theta}) - \mathcal{V}(t'\hat{\theta})] = t'[\mathcal{V}(\tilde{\theta}) - \mathcal{V}(\hat{\theta})]t$. □

Example 11.3 With $2N$ i.i.d. observations $y_i \sim \mathcal{N}(\mu, \sigma^2)$, the estimators

$$\tilde{\mu}_1 = \sum_{i=1}^N y_i / N \tag{11.51}$$

$$\tilde{\mu}_2 = \sum_{i=N+1}^{2N} y_i / N \tag{11.52}$$

$$\hat{\mu} = \sum_{i=1}^{2N} y_i / (2N) = (\tilde{\mu}_1 + \tilde{\mu}_2) / 2 \tag{11.53}$$

are all unbiased. However, $\mathcal{V}(\tilde{\mu}_1) = \mathcal{V}(\tilde{\mu}_2) = 2\mathcal{V}(\hat{\mu})$.

Corollary 11.14 For a $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$ with $r = \text{rank}(X) \leq q$, the BLUE of estimable $\theta = C\beta (a \times 1)$ is $\hat{\theta} \Leftrightarrow$ the BLUE of $t'\theta$ is $t'\hat{\theta} \forall t \in \mathbb{R}^a$.

Proof. The BLUE properties linearity, $\hat{\theta} = T\mathbf{y} \Leftrightarrow t'\hat{\theta} = t'T\mathbf{y} \forall t \in \mathbb{R}^a$, and unbiasedness give $E(\hat{\theta}) = \theta \Leftrightarrow t'E(\hat{\theta}) = t'\theta \forall t \in \mathbb{R}^a$. To prove the variance is a minimum, we proceed as follows. The matrix $\mathcal{V}(\hat{\theta})$ is minimal (as in BLUE version a) $\Leftrightarrow \forall t \in \mathbb{R}^a$ scalar $\mathcal{V}(t'\hat{\theta})$ is minimal. □

Theorem 11.15 For a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with $r = \text{rank}(\mathbf{X}) \leq q$, the unique BLUE of an estimable $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ ($a \times 1$) is $\hat{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}}$, with $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$ the least squares estimator of $\boldsymbol{\beta}$.

The assumptions guarantee $\hat{\boldsymbol{\theta}}$ is unbiased, is invariant to the choice of generalized inverse defining $\tilde{\boldsymbol{\beta}}$, has minimum variance among all LUEs, and is unique in that it is the only estimator satisfying all of the requirements.

Proof. It is sufficient to prove $\forall t \in \Re^a$ the BLUE of $t'\boldsymbol{\theta}$ is $t'\hat{\boldsymbol{\theta}}$. If $\mathbf{c}' = t'\mathbf{C}$, then by Lemma 11.7 the BLUE of $\theta_t = \mathbf{c}'\boldsymbol{\beta} = t'\mathbf{C}\boldsymbol{\beta} = t'\boldsymbol{\theta}$ is $\hat{\theta}_t = (t'\mathbf{C})(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$. \square

Corollary 11.15 For the special case of a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with estimable $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ the unique BLUE is $\hat{\boldsymbol{\theta}} = \mathbf{C}\hat{\boldsymbol{\beta}}$ in which $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$.

11.11 RELATED DISTRIBUTIONS

It is definitely worth repeating that nearly all results about estimation presented earlier in the present chapter, except for likelihood properties, are distribution free. As long as the GLM assumptions hold, the random variables can have any distribution with finite second moments. The special case of Gaussian errors leads to simple forms of distributions for many estimators.

Theorem 11.16 (a) For $\text{GLM}_{N,q}\text{LTFR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors and two-condition inverse $(\mathbf{X}'\mathbf{X})^{-}$,

$$\tilde{\boldsymbol{\beta}} \sim \text{SN}_q[(\mathbf{X}'\mathbf{X})^{-}(\mathbf{X}'\mathbf{X})\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-}]. \quad (11.54)$$

(b) A $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ has

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}_q[\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]. \quad (11.55)$$

(c) For $\text{rank}(\mathbf{X}) \leq q$ and any one-condition inverse, estimable $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ ensures

$$\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 \sim (S)\mathcal{N}_a[\boldsymbol{\theta} - \boldsymbol{\theta}_0, \sigma^2\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}']. \quad (11.56)$$

(d) For any one-condition inverse, $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' = \mathbf{X}(\mathbf{X}'\mathbf{X})^{+}\mathbf{X}'$ and

$$\begin{aligned} \hat{\mathbf{y}} &= \mathbf{X}\tilde{\boldsymbol{\beta}} = \mathbf{H}\mathbf{y} \\ &\sim \text{SN}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{H}). \end{aligned} \quad (11.57)$$

(e) Furthermore

$$\begin{aligned} \hat{\mathbf{e}} = \mathbf{y} - \hat{\mathbf{y}} &= (\mathbf{I}_N - \mathbf{H})\mathbf{y} \\ &\sim \text{SN}_N[\mathbf{X}\boldsymbol{\beta}, \sigma^2(\mathbf{I}_N - \mathbf{H})]. \end{aligned} \quad (11.58)$$

Proof. Left as an exercise. *Hints:* linear transformations of Gaussian variables, Theorem 1.15, generalized inverse conditions, and estimability condition.

Theorem 11.17 In a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors, $\tilde{\boldsymbol{\beta}}$ and $\hat{\boldsymbol{e}}$ are jointly Gaussian with zero covariance and correlation and are statistically independent.

Proof. The joint distribution of the estimators is established by expressing the two vectors as a single linear transformation of a vector Gaussian:

$$\begin{bmatrix} \tilde{\boldsymbol{\beta}} \\ \hat{\boldsymbol{e}} \end{bmatrix} = \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}' \\ (\mathbf{I} - \mathbf{H}) \end{bmatrix} \mathbf{y} = \mathbf{T}\mathbf{y}. \tag{11.59}$$

In turn

$$\begin{aligned} \mathcal{V}(\tilde{\boldsymbol{\beta}}, \hat{\boldsymbol{e}}) &= \text{E}(\tilde{\boldsymbol{\beta}}\hat{\boldsymbol{e}}') - \text{E}(\tilde{\boldsymbol{\beta}})\text{E}(\hat{\boldsymbol{e}}') \\ &= \text{E}\{(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}[(\mathbf{I} - \mathbf{H})\mathbf{y}]'\} - \text{E}(\tilde{\boldsymbol{\beta}})\mathbf{0} \\ &= (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\text{E}(\mathbf{y}\mathbf{y}')(\mathbf{I} - \mathbf{H}) \\ &= (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'(\sigma^2\mathbf{I}_N)(\mathbf{I} - \mathbf{H}) \\ &= \sigma^2(\mathbf{X}'\mathbf{X})^{-}[\mathbf{X}' - \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'] = \sigma^2(\mathbf{X}'\mathbf{X})^{-}[\mathbf{0}]. \end{aligned} \tag{11.60}$$

The last step is true by Theorem 1.15 and gives $\mathcal{V}(\tilde{\boldsymbol{\beta}}, \hat{\boldsymbol{e}}) = \mathbf{0}$. □

11.12 FORMULATIONS OF EXPLICIT RESTRICTIONS OF $\boldsymbol{\beta}$ AND $\boldsymbol{\theta}$

The notation $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ describes a situation in which no restrictions have been placed on relationships among the elements of $\boldsymbol{\beta}$. In some cases, data analysts wish to ensure some explicit linear relationships hold among the elements of $\boldsymbol{\beta}$. The classical approach to one-way ANOVA coding for G groups uses the model $y_{ij} = \mu + \alpha_j + e_{ij}$ subject to the restriction $\sum_{j=1}^G \alpha_j = 0$. The use of the restriction overcomes the inherently LTFR nature of the classical coding. It leads to a method for identifying a set of G estimable parameters (from among the set of $G + 1$ nonestimable original parameters).

The notation may be extended to incorporate explicit linear restrictions on the elements of $\boldsymbol{\beta}$ by writing $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, |\mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$ for \mathbf{R} ($k \times q$) and \mathbf{a} ($k \times 1$) conforming and known constants with $k \leq q$. Most often, $\mathbf{a} = \mathbf{0}$. Without restrictions, $\boldsymbol{\beta} \in \mathfrak{R}^q$, whereas, with the explicit restrictions, $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$, $\boldsymbol{\beta}$ is required to lie in a subspace. Such linear restrictions on (the parameters space for) $\boldsymbol{\beta}$ can be imposed on both FR and LTFR models. They are particularly important for LTFR models because such additional restrictions can lead to a model with *estimable* primary expected-value parameters. They will also be seen to be important in the formulation of the general linear hypothesis test.

Example 11.4 The classical ANOVA coding for a $G = 3$ group one-way ANOVA provides a simple example. In matrix notation, the model is

$$\begin{aligned}
 \mathbf{y} &= \mathbf{X}\boldsymbol{\beta} + \mathbf{e} \\
 &= \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} + \mathbf{e}.
 \end{aligned}
 \tag{11.61}$$

The “sum to zero” zero constraints may be stated as

$$\begin{aligned}
 \mathbf{R}\boldsymbol{\beta} &= \mathbf{a} \\
 \begin{bmatrix} 0 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \\ \alpha_3 \end{bmatrix} &= [0].
 \end{aligned}
 \tag{11.62}$$

Since the usual estimation procedures involve optimization, estimation subject to linear restrictions can be performed via constrained optimization, such as through the use of Lagrange multipliers. However, the following two theorems allow restating the restrictions and in turn restating the restricted model as an unrestricted model.

Theorem 11.18 A restricted GLM_{N,q}($y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2$) has \mathbf{R} ($k \times q$) and \mathbf{a} ($k \times 1$) known constants, $k \leq q$, and consistent equations $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$ (with at least one solution).

(a) If $\text{rank}(\mathbf{R}) \leq k$ and \mathbf{R}^- ($q \times k$) is any particular generalized inverse (i.e., $\mathbf{R}\mathbf{R}^-\mathbf{R} = \mathbf{R}$), then any value $\boldsymbol{\beta}_R$ which satisfies the restrictions may be written, for some $\boldsymbol{\gamma} \in \mathbb{R}^q$, as

$$\boldsymbol{\beta}_R = \mathbf{R}^-\mathbf{a} + (\mathbf{R}^-\mathbf{R} - \mathbf{I}_q)\boldsymbol{\gamma}.
 \tag{11.63}$$

(b) If $\text{rank}(\mathbf{R}) = k$, then the singular value decomposition allows writing

$$\begin{aligned}
 \mathbf{R} &= \mathbf{P}[\boldsymbol{\Lambda} \mathbf{0}]\mathbf{Q}' \\
 &= \mathbf{P}[\boldsymbol{\Lambda} \mathbf{0}][\mathbf{Q}_1 \mathbf{Q}_2] \\
 &= \mathbf{P}\boldsymbol{\Lambda}\mathbf{Q}'_1,
 \end{aligned}
 \tag{11.64}$$

the (unique) Moore-Penrose (four-condition) generalized inverse is

$$\mathbf{R}^+ = \mathbf{Q}_1\boldsymbol{\Lambda}^{-1}\mathbf{P}',
 \tag{11.65}$$

and $\boldsymbol{\beta}$ satisfies the restrictions if and only if, for some $\boldsymbol{\tau} \in \mathbb{R}^{q-k}$,

$$\boldsymbol{\beta} = \mathbf{R}^+\mathbf{a} + \mathbf{Q}_2\boldsymbol{\tau}.
 \tag{11.66}$$

Here \mathbf{P} is $k \times k$, \mathbf{Q} is $q \times q$, \mathbf{Q}_1 is $q \times k$, and \mathbf{Q}_2 with $q \times (q - k)$, with all being columnwise orthonormal.

(c) If $\text{rank}(\mathbf{R}) = k$, then, with $s = q - k$ and columns of \mathbf{R} and corresponding rows of $\boldsymbol{\beta}$ permuted prior to partitioning if necessary,

$$\mathbf{R}\boldsymbol{\beta} = [\mathbf{R}_1 \ \mathbf{R}_2] \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix}, \tag{11.67}$$

with \mathbf{R}_1 $k \times k$ and full rank, which ensures \mathbf{R}_1^{-1} exists. Also $\boldsymbol{\beta}_1$ is $k \times 1$, \mathbf{R}_2 is $k \times s$, and $\boldsymbol{\beta}_2$ is $s \times 1$. Only $\boldsymbol{\beta}_2$ is free to vary and $\boldsymbol{\beta}$ satisfies $\mathbf{R}\boldsymbol{\beta} = \mathbf{a} \Leftrightarrow$

$$\boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{R}_1^{-1} \\ \mathbf{0}_{s \times k} \end{bmatrix} \mathbf{a} + \begin{bmatrix} -\mathbf{R}_1^{-1} \mathbf{R}_2 \\ \mathbf{I}_s \end{bmatrix} \boldsymbol{\beta}_2. \tag{11.68}$$

Proof of (a). (\Leftarrow) Any such $\boldsymbol{\beta}_R$ is a solution.

(\Rightarrow) If $\boldsymbol{\beta}$ is any solution, then choosing $\boldsymbol{\gamma} = (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)\boldsymbol{\beta}$ ensures $\boldsymbol{\beta}$ satisfies the relationship required of $\boldsymbol{\beta}_R$. In turn,

$$\begin{aligned} \mathbf{R}^- \mathbf{a} + (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)\boldsymbol{\gamma} &= \mathbf{R}^- \mathbf{a} + (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)\boldsymbol{\beta} \\ &= \mathbf{R}^- \mathbf{a} + (\mathbf{R}^- \mathbf{R} \mathbf{R}^- \mathbf{R} - \mathbf{R}^- \mathbf{R} - \mathbf{R}^- \mathbf{R} + \mathbf{I}_q)\boldsymbol{\beta} \\ &= \mathbf{R}^- \mathbf{a} + (\mathbf{I}_q - \mathbf{R}^- \mathbf{R})\boldsymbol{\beta} \\ &= \mathbf{R}^- (\mathbf{R}\boldsymbol{\beta}) + (\mathbf{I}_q - \mathbf{R}^- \mathbf{R})\boldsymbol{\beta} \\ &= \boldsymbol{\beta} + \mathbf{0}. \end{aligned} \tag{11.69}$$

Proof of (b). If $\mathbf{R}^- = \mathbf{R}^+$ and $\boldsymbol{\tau} = \mathbf{Q}'_2 \boldsymbol{\gamma}$, applying part (a) gives

$$\begin{aligned} \mathbf{R}^- \mathbf{a} + (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)\boldsymbol{\gamma} &= \mathbf{R}^+ \mathbf{a} + (\mathbf{Q}_1 \mathbf{Q}'_1 - \mathbf{I}_q)\boldsymbol{\gamma} \\ &= \mathbf{R}^+ \mathbf{a} + (\mathbf{Q}_2 \mathbf{Q}'_2)\boldsymbol{\gamma} \\ &= \mathbf{R}^+ \mathbf{a} + \mathbf{Q}_2(\boldsymbol{\tau}). \end{aligned} \tag{11.70}$$

Proof of (c). $\mathbf{R}_1 \boldsymbol{\beta}_1 + \mathbf{R}_2 \boldsymbol{\beta}_2 = \mathbf{a}$ implies $\boldsymbol{\beta}_1 = \mathbf{R}_1^{-1} \mathbf{a} - \mathbf{R}_1^{-1} \mathbf{R}_2 \boldsymbol{\beta}_2$. □

Corollary 11.18 The matrix $[\mathbf{X}' \ \mathbf{R}']$ has full rank $\Leftrightarrow \mathbf{Z} = (\mathbf{X}_2 - \mathbf{X}_1 \mathbf{R}_1^{-1} \mathbf{R}_2)$ has full rank.

Proof. Left as an exercise.

11.13 RESTRICTED ESTIMATION VIA EQUIVALENT MODELS

Definition 11.10 (a) An *explicitly restricted linear model* is indicated $\text{GLM}_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$, with $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$ indicating restrictions on $\boldsymbol{\beta}$.
(b) A linearly equivalent unrestricted form is an *implicitly restricted linear model*.

The following theorem provides explicit methods for finding an unrestricted model for any given restricted model. Hence the theorem guarantees a linearly equivalent unrestricted form always exists for any restricted model.

Theorem 11.19 A restricted GLM $_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$ has \mathbf{R} ($k \times q$) and \mathbf{a} ($k \times 1$) known constants, with $k \leq q$, $r = \text{rank}(\mathbf{X}) \leq q$ and $s = q - k$. Consistent equations $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$ give three distinct ways to create and analyze a linearly equivalent unrestricted model, depending on $\text{rank}(\mathbf{R})$.

(a) Having $\text{rank}(\mathbf{R}) \leq k$ allows the following method.

1. Compute \mathbf{R}^- ($q \times k$) such that $\mathbf{R}\mathbf{R}^- \mathbf{R} = \mathbf{R}$.
2. Compute transformed responses $\mathbf{u} = \mathbf{y} - \mathbf{X}\mathbf{R}^- \mathbf{a}$ ($N \times 1$).
3. Compute transformed predictors $\mathbf{Z} = \mathbf{X}(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)$ ($N \times q$).
4. Analyze unrestricted model (I) GLM $_{N,q}(u_i; \mathbf{Z}_i\boldsymbol{\gamma}, \sigma^2)$ with $\text{rank}(\mathbf{Z}) \leq s$.

(b) Having $\text{rank}(\mathbf{R}) = k$ allows the following method.

1. Compute the SVD $\mathbf{R} = \mathbf{P}[\boldsymbol{\Lambda} \mathbf{0}][\mathbf{Q}_1 \mathbf{Q}_2]'$ with $\mathbf{P}'\mathbf{P} = \mathbf{P}\mathbf{P}' = \mathbf{I}_k$, \mathbf{Q}_1 ($q \times k$), $\mathbf{Q}_1'\mathbf{Q}_1 = \mathbf{I}_k$, and \mathbf{Q}_2 ($q \times s$), $\mathbf{Q}_2'\mathbf{Q}_2 = \mathbf{I}_s$.
2. Compute $\mathbf{R}^+ = \mathbf{Q}_1\boldsymbol{\Lambda}^{-1}\mathbf{P}'$.
3. Compute transformed responses $\mathbf{u} = \mathbf{y} - \mathbf{X}\mathbf{R}^+ \mathbf{a}$ ($N \times 1$).
4. Compute transformed (and reduced) predictors $\mathbf{Z} = \mathbf{X}\mathbf{Q}_2$ ($N \times s$).
5. Analyze unrestricted model (II) GLM $_{N,s}(u_i; \mathbf{Z}_i\boldsymbol{\tau}, \sigma^2)$ with $\text{rank}(\mathbf{Z}) \leq s$.

(c) Having $\text{rank}(\mathbf{R}) = k$ also allows the following method.

1. Permute columns of \mathbf{X} and corresponding rows of $\boldsymbol{\beta}$ to have $\mathbf{R} = [\mathbf{R}_1 \mathbf{R}_2]$, with \mathbf{R}_1 ($k \times k$) full rank and $\mathbf{X} = [\mathbf{X}_1 \mathbf{X}_2]$ partitioned similarly.
2. Compute transformed responses $\mathbf{u} = \mathbf{y} - \mathbf{X}_1\mathbf{R}_1^{-1}\mathbf{a}$ ($N \times 1$).
3. Compute transformed (and reduced) predictors ($N \times s$)

$$\mathbf{Z} = \mathbf{X} \begin{bmatrix} -\mathbf{R}_1^{-1}\mathbf{R}_2 \\ \mathbf{I}_s \end{bmatrix}. \tag{11.71}$$

4. Analyze unrestricted model (III) GLM $_{N,s}(u_i; \mathbf{Z}_i\boldsymbol{\beta}_2, \sigma^2)$ with $\text{rank}(\mathbf{Z}) \leq s$. Here $(\mathbf{u} - \mathbf{Z}\boldsymbol{\gamma}) = (\mathbf{u} - \mathbf{Z}\boldsymbol{\tau}) = (\mathbf{u} - \mathbf{Z}\boldsymbol{\beta}_2) = (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})$.

Proof. If $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon}$ given $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$, then Theorem 11.18 (a) implies

$$\mathbf{y} = \mathbf{X}\mathbf{R}^- \mathbf{a} + \mathbf{X}(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)\boldsymbol{\gamma} + \boldsymbol{\epsilon}. \tag{11.72}$$

Hence the restricted model and model I are linearly equivalent. Theorem 11.18 (b) implies

$$\mathbf{y} = \mathbf{X}\mathbf{R}^+ \mathbf{a} + \mathbf{X}\mathbf{Q}_2\boldsymbol{\tau} + \boldsymbol{\epsilon}. \tag{11.73}$$

Hence the restricted model and model II are linearly equivalent. Theorem 11.18 (c) implies

$$\mathbf{y} = \mathbf{X}\mathbf{R}_1^{-1}\mathbf{a} + \mathbf{X} \begin{bmatrix} -\mathbf{R}_1^{-1}\mathbf{R}_2 \\ \mathbf{I}_s \end{bmatrix} \boldsymbol{\beta}_2 + \boldsymbol{\epsilon}. \tag{11.74}$$

Hence the restricted model and model III are linearly equivalent. □

Theorem 11.20 With the notation and assumptions of Theorem 11.18 (b), $\theta = C\beta$ ($a \times 1$) is estimable. Here $C = I$ is allowed $\Leftrightarrow \text{rank}(X) = q$. The BLUE of θ is $\hat{\theta}_R = C(R^+a + Q_2\bar{\tau})$, in which $\bar{\tau}$ is the BLUE of τ ($s \times 1$) in the linearly equivalent unrestricted model II, namely $\text{GLM}_{N,s}(u_i; Z_i\tau, \sigma^2)$ with $\text{rank}(Z) \leq s$, $u = y - XR^+a$, and $Z = XQ_2$ ($N \times s$).

Proof. Subject to the restrictions $R\beta = a$, $\theta = C\beta$ is still estimable because $\theta = CR^+a + CQ_2\tau$ implies $\theta = \text{constant} + (CQ_2)\tau$. Also, $CQ_2\tau$ is estimable in unrestricted model II because $C\beta$ is estimable in the restricted model. Therefore, by Theorem 11.15, the BLUE of $CQ_2\tau$ is $CQ_2\bar{\tau}$, and the desired result follows. \square

Corollary 11.20 (a) If $\text{rank}(X) = q$, then the BLUE of β is

$$\begin{aligned} \hat{\beta}_R &= (R^+a + Q_2\bar{\tau}) \\ &= \hat{\beta} - (X'X)^{-1}R' \left[R(X'X)^{-1}R' \right]^{-1} (R\hat{\beta} - a), \end{aligned} \tag{11.75}$$

in which $\hat{\beta} = (X'X)^{-1}X'y$.

(b) An unbiased estimator of σ^2 is

$$\begin{aligned} \hat{\sigma}^2 &= (y - X\hat{\beta}_R)'(y - X\hat{\beta}_R) / [(N - q + k)] \\ &= (SSE + SSH) / [(N - q + k)], \end{aligned} \tag{11.76}$$

with

$$SSE = y'[I - X(X'X)^{-1}X']y = (y - X\hat{\beta})'(y - X\hat{\beta}) \tag{11.77}$$

and

$$SSH = (R\hat{\beta} - a)' \left[R(X'X)^{-1}R' \right]^{-1} (R\hat{\beta} - a). \tag{11.78}$$

Proof. The first form for $\hat{\beta}_R$ may be found by considering $\theta = \beta$ and applying the original theorem. In the special case of $\text{rank}(X) = q$ the parameter $\theta = C\beta$ is estimable for all C including $C = I$.

The second form for $\hat{\beta}_R$ may be found by using Theorem 11.18 (a) to help define a linearly equivalent unrestricted model, $\text{GLM}_{N,q}(u_i; Z_i\gamma, \sigma^2)$ with Z $N \times q$ with $\text{rank}(Z) \leq q$, R $q \times k$, a $k \times 1$, $u = y - XR^-a$, $Z = X(R^-R - I_q)$, and $\beta_R = R^-a + (R^-R - I_q)\gamma$. Here β_R is an estimable parameter, $\beta_R = C\gamma + R^-a$, since $C = (R^-R - I_q) = (X'X)^{-1}X'Z$ is a linear combination of the rows of Z , and R^-a is a constant. By Theorem 11.15 the unique BLUE of β_R is $\hat{\beta}_R = C\tilde{\gamma} + R^-a$ in which $\tilde{\gamma} = (Z'Z)^-Z'U$. The estimator can be written

$$\begin{aligned} \widehat{\beta}_R &= (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'(\mathbf{Y} - \mathbf{X}\mathbf{R}^- \mathbf{a}) + \mathbf{R}^- \mathbf{a} \\ &= [(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}']\mathbf{y} + \{[\mathbf{I}_k - (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}'\mathbf{X}]\mathbf{R}^- \} \mathbf{a} \\ &= \mathbf{T}_{1y}\mathbf{y} + \mathbf{T}_{1a}\mathbf{a}. \end{aligned} \tag{11.79}$$

The estimator can also be written

$$\begin{aligned} \widehat{\beta}_R &= \widehat{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}(\mathbf{R}\widehat{\beta} - \mathbf{a}) \\ &= \{(\mathbf{I}_q - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\mathbf{R})\}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} + \\ &\quad \{(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\}\mathbf{a} \\ &= \mathbf{T}_{2y}\mathbf{y} + \mathbf{T}_{2a}\mathbf{a}. \end{aligned} \tag{11.80}$$

The equality of the expressions can be proved by first equating the coefficients of \mathbf{y} and second equating the coefficients of \mathbf{a} .

The fact that \mathbf{R}^- can be chosen to ensure $(\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}' = (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ helps accomplish the two equating tasks. The forms

$$\begin{aligned} \mathbf{Z}^- &= (\mathbf{Z}'\mathbf{Z})^{-1}\mathbf{Z}' && \text{(i)} \\ \mathbf{Z}^- &= (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' && \text{(ii)} \end{aligned}$$

are both valid representations for a generalized inverse of \mathbf{Z} . The representation (i) is contained in Theorem 1.15. The validity of (ii) follows directly from verifying $\mathbf{Z}\mathbf{Z}^-\mathbf{Z} = \mathbf{Z}$ for $\mathbf{Z} = \mathbf{X}(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)$. All possible generalized inverses of \mathbf{Z} are represented by each of the two expressions. Therefore, for any particular generalized inverse of $\mathbf{Z}'\mathbf{Z}$ a particular generalized inverse of \mathbf{R} exists which makes (i) and (ii) identical.

Applying the equality, we have

$$\begin{aligned} \mathbf{T}_{1y} &= (\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{R}^- \mathbf{R} - \mathbf{I}_q)(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ &= (\mathbf{I}_q - \mathbf{R}^- \mathbf{R})(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}', \end{aligned} \tag{11.81}$$

compared with $\mathbf{T}_{2y} = \{ \mathbf{I}_q - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\mathbf{R} \}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$. Hence it will suffice to prove $\mathbf{R}^- \mathbf{R} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\mathbf{R}$. Here

$$\begin{aligned} \mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}' &= \mathbf{R}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'] && \downarrow \\ \mathbf{R}^{-1}[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'] &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}' && \downarrow \\ \mathbf{R}^- &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1} && \downarrow \\ \mathbf{R}^- \mathbf{R} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}'[\mathbf{R}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{R}']^{-1}\mathbf{R} && \downarrow \\ \mathbf{T}_{1y} &= \mathbf{T}_{2y}. \end{aligned} \tag{11.82}$$

Given $\mathbf{T}_{1y} = \mathbf{T}_{2y}$, it follows almost immediately that $\mathbf{T}_{1a} = \mathbf{T}_{2a}$ since

$$\begin{aligned} \mathbf{T}_{1a} &= (\mathbf{I} - \mathbf{T}_{1y}\mathbf{X})\mathbf{R}^- \\ \mathbf{T}_{2a}\mathbf{R} &= \mathbf{I} - \mathbf{T}_{2y}\mathbf{X}. \end{aligned} \tag{11.83}$$

Multiplying the first equation by \mathbf{a} and the second equation by $\mathbf{R}^- \mathbf{a}$ gives

$$\begin{aligned} T_{1a}a &= (I - T_{1y}X)R^{-}a \\ T_{2a}RR^{-}a &= (I - T_{2y}X)R^{-}a. \end{aligned} \tag{11.84}$$

Since $RR^{-}a = a$ and $T_{1y} = T_{2y}$, it follows that $T_{1a}a = T_{2a}a$ for all a . Thus $T_{1y} = T_{2y}$ and $T_{1a} = T_{2a}$.

It is not necessary to assume Gaussian errors to formulate a reasonable estimator of σ^2 . By Theorem 11.18, a linearly equivalent unrestricted model is $GLM_{N,q}(u_i; Z_i\gamma, \sigma^2)$ with $\text{rank}(Z) \leq q$, $u = y - XR^{-}a$, $Z = X(R^{-}R - I_q)$, and $\beta_R = R^{-}a + (R^{-}R - I_q)\gamma$. For any real symmetric A $E(u' Au) = \text{tr}(A\sigma^2 I_N) + (Z\tau)' A (Z\tau)$. If $A = [I - Z(Z'Z)^{-}Z] \equiv (I - P_z)$, then $E(u' Au) = \sigma^2 \text{rank}(A) + 0 = \sigma^2[N - (q - k)]$. Hence an unbiased estimator for σ^2 is $\hat{\sigma}^2 = u'(I - P_z)u / (N - q + k) = (u - Z\hat{\tau})'(u - Z\hat{\tau}) / (N - q + r)$. As noted in Theorem 11.18, $(u - Z\hat{\tau}) = (y - X\hat{\beta}_R)$ and hence the desired result follows.

The second form for $\hat{\sigma}^2$ may be found in terms of the following expressions:

$$(\hat{\beta} - \hat{\beta}_R) = (X'X)^{-}R'[R(X'X)^{-}R']^{-}(R\hat{\beta} - a) \tag{11.85}$$

$$SSE = (y - X\hat{\beta})'(y - X\hat{\beta}) \tag{11.86}$$

$$SSH = (R\hat{\beta} - a)' [R(X'X)^{-}R']^{-}(R\hat{\beta} - a). \tag{11.87}$$

Simple manipulations give

$$\begin{aligned} \hat{\sigma}^2(N - q + k) &= (y - X\hat{\beta}_R)'(y - X\hat{\beta}_R) \\ &= [y - X\hat{\beta} + X(\hat{\beta} - \hat{\beta}_R)]' [y - X\hat{\beta} + X(\hat{\beta} - \hat{\beta}_R)] \\ &= (y - X\hat{\beta})'(y - X\hat{\beta}) + (\hat{\beta} - \hat{\beta}_R)'(X'X)(\hat{\beta} - \hat{\beta}_R) \\ &= SSE + SSH. \end{aligned} \tag{11.88}$$

□

11.14 FITTING PIECEWISE POLYNOMIAL MODELS VIA SPLINES

Definition 11.11 A collection of m piecewise linear functions connected at m points, known as *knots*, together define a *spline* function.

In some cases, the study design and scientific setting dictate that no single model applies to the entire range of the predictor. The complication could arise in modeling the yield of a chemical manufacturing process as a function of total volume of ingredients. The batch size z_i (a model predictor) can vary only within a limited range for each of $m = 3$ container sizes. Using regression splines allows specifying a distinct polynomial model for each container size, as does the obvious alternative of fitting three different models. However, not only do splines allow creating a smooth model for the entire range of interest, they also provide much greater precision by using a pooled estimator of variance. Smith (1979) provided a useful introduction to regression splines.

A statistical model for \mathbf{y} ($N \times 1$) specifying the mean as a piecewise polynomial function with $m = 3$ parts provides an example:

$$E(y_i|z_i) = f(z_i; \boldsymbol{\beta}) = \begin{cases} a_1 + b_1 z_i & z_i \in (k_0, k_1] \\ a_2 + b_2 z_i + c_2 z_i^2 & z_i \in (k_1, k_2] \\ a_3 + b_3 z_i & z_i \in (k_2, k_3]. \end{cases} \quad (11.89)$$

The points of connection \mathbf{k} ($m \times 1$) are the knots. The m values are assumed to be fixed and known. Usually the values are chosen by the analyst and are evenly spaced with $m \leq N^{1/5}$. With the knots fixed and known, the only remaining parameters of $f(\cdot)$ are the regression coefficients $\boldsymbol{\beta} = [a_1 \ b_1 \ a_2 \ b_2 \ c_2 \ a_3 \ b_3]'$. The fact that $E(y_i|z_i) = f(z_i; \boldsymbol{\beta}) = \mathbf{X}_i \boldsymbol{\beta}$ verifies the mean is linear in $\boldsymbol{\beta}$. If $\mathbf{X}_{1i} = [1 \ z_i]$, $\delta_{1i} = \mathbb{I}\{z_i \in (k_0, k_1]\}$, $\mathbf{X}_{2i} = [1 \ z_i \ z_i^2]$, $\delta_{2i} = \mathbb{I}\{z_i \in (k_1, k_2]\}$, $\mathbf{X}_{3i} = [1 \ z_i]$, $\delta_{3i} = \mathbb{I}\{z_i \in (k_2, k_3]\}$, $\boldsymbol{\beta}_1 = [a_1 \ b_1]'$, $\boldsymbol{\beta}_2 = [a_2 \ b_2 \ c_2]'$, and $\boldsymbol{\beta}_3 = [a_3 \ b_3]'$, then equivalent expressions for $f(z_i; \boldsymbol{\beta})$ are

$$\begin{aligned} f(z_i; \boldsymbol{\beta}) &= \begin{cases} \mathbf{X}_{1i} \boldsymbol{\beta}_1 & z_i \in (k_0, k_1] \\ \mathbf{X}_{2i} \boldsymbol{\beta}_2 & z_i \in (k_1, k_2] \\ \mathbf{X}_{3i} \boldsymbol{\beta}_3 & z_i \in (k_2, k_3] \end{cases} \\ &= \sum_{j=1}^3 \delta_{ji} \mathbf{x}_{ji} \boldsymbol{\beta}_j \\ &= [\delta_{1i} \mathbf{X}_{1i} \quad \delta_{2i} \mathbf{X}_{2i} \quad \delta_{3i} \mathbf{X}_{3i}] \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \boldsymbol{\beta}_3 \end{bmatrix} \\ &= \mathbf{X}_i \boldsymbol{\beta}. \end{aligned} \quad (11.90)$$

It follows $E(\mathbf{y}|\mathbf{z}) = f(\mathbf{z}; \boldsymbol{\beta}) = \mathbf{X} \boldsymbol{\beta}$. (Can you say how \mathbf{X} is defined?)

Theorem 11.21 Regression splines allow defining a valid univariate GLM. We assume $E(y_i|z_i)$ is a piecewise linear function of $\boldsymbol{\beta}$, namely

$$E(y_i|z_i) = \begin{cases} \mathbf{X}_{i,1} \boldsymbol{\beta}_1 & z_i \in (k_0, k_1] \\ \mathbf{X}_{i,2} \boldsymbol{\beta}_2 & z_i \in (k_1, k_2] \\ \vdots & \\ \mathbf{X}_{i,m} \boldsymbol{\beta}_m & z_i \in (k_{m-1}, k_m] \end{cases} \quad (11.91)$$

with known points of connection (knots) \mathbf{k} ($m \times 1$). If

$$d_{i,j} = \begin{cases} 1 & z_i \in (k_{j-1}, k_j] \\ 0 & \text{otherwise,} \end{cases} \quad (11.92)$$

$$\mathbf{X}_i = [d_{i1} \mathbf{X}_{i,1} \ \cdots \ d_{im} \mathbf{X}_{i,m}] \quad (11.93)$$

and

$$\boldsymbol{\beta} = \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_m \end{bmatrix}, \tag{11.94}$$

then \mathbf{y} satisfies $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$, with q the number of columns in $\{\mathbf{X}_j\}$.

Proof. Left as an exercise.

Generally, the polynomial pieces $\{f_j(z_i; \boldsymbol{\beta}), j \in \{1, \dots, m\}\}$ do not connect and will not form a smooth continuous curve unless some constraints are placed on them. To make $f(z_i; \boldsymbol{\beta}) = \sum_{j=1}^m \delta_{ji} f_j(z_i; \boldsymbol{\beta})$ continuous at knot k_j we must require $f_j(z; \boldsymbol{\beta})|_{z=k_j} = f_{j+1}(z; \boldsymbol{\beta})|_{z=k_j}$. Even with the requirement, the pieces may not connect to form a smooth curve. A degree of smoothness is obtained by requiring the derivatives from the left and right to be equal,

$$\left. \frac{\partial f_j(z; \boldsymbol{\beta})}{\partial z} \right|_{z=k_j} = \left. \frac{\partial f_{j+1}(z; \boldsymbol{\beta})}{\partial z} \right|_{z=k_j}. \tag{11.95}$$

Higher degrees of smoothness are obtained by specifying requirements for higher order derivatives such as

$$\left. \frac{\partial^2 f_j(z; \boldsymbol{\beta})}{\partial z^2} \right|_{z=k_j} = \left. \frac{\partial^2 f_{j+1}(z; \boldsymbol{\beta})}{\partial z^2} \right|_{z=k_j}. \tag{11.96}$$

Each constraint at each knot constitutes a single linear constraint on $\boldsymbol{\beta}$. Collectively the constraints can be expressed in the usual form, $\mathbf{R}\boldsymbol{\beta} = \mathbf{a}$. The resulting constrained model can be transformed to a linearly equivalent unconstrained model.

Corollary 11.21 If smoothness constraints are placed on $f(z_i; \boldsymbol{\beta})$, then \mathbf{y} satisfies $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \sigma^2)$.

Proof. Left as an exercise.

Some scientific settings require estimating the join points rather than having them fixed. Gallant and Fuller (1973) provided a useful discussion of splines.

11.15 ESTIMATION FOR THE GGLM: WEIGHTED LEAST SQUARES

All results presented so far in the chapter apply not only to the GLM but also to the GGLM. The generality arises due to most of the results depending only on properties of the design matrix as they relate to expected values.

Definition 11.12 Exactly as for two GLMs, models $\text{GGLM}_{N,q_1}(\mathbf{y}; \mathbf{X}_1\boldsymbol{\beta}_1, \boldsymbol{\Upsilon})$ and $\text{GGLM}_{N,q_2}(\mathbf{y}; \mathbf{X}_2\boldsymbol{\beta}_2, \boldsymbol{\Upsilon})$ are *linearly equivalent* whenever (1) for any $\boldsymbol{\beta}_1$ there exists $\boldsymbol{\beta}_2$ such that $\mathbf{X}_1\boldsymbol{\beta}_1 = \mathbf{X}_2\boldsymbol{\beta}_2$ and (2) for any $\boldsymbol{\beta}_2$ there exists $\boldsymbol{\beta}_1$ such that $\mathbf{X}_1\boldsymbol{\beta}_1 = \mathbf{X}_2\boldsymbol{\beta}_2$.

It is useful to recognize that a GGLM may be linearly equivalent to a univariate GLM. Theorems about weighted least squares later in the chapter use such an equivalence. Similar results lie at the heart of some derivations of the “univariate” approach to repeated measures. The approach depends on the assumption of compound symmetric covariance among a set of p observations on each of N independent sampling units.

Theorem 11.22 For any $\text{GGLM}_{N,q}\text{LTFR}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \boldsymbol{\Upsilon})$ with $r = \text{rank}(\mathbf{X})$, linearly equivalent $\text{GGLM}_{N,r}\text{FR}(\mathbf{y}; \mathbf{X}_1\boldsymbol{\beta}_1, \boldsymbol{\Upsilon})$ always exists.

Proof. Left as an exercise.

Although results about restricted models are stated in terms of the univariate GLM, they also apply to the $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{a}, \boldsymbol{\Upsilon})$. The generalization is valid because the constraints apply only to the expected-value portion of the model.

Definition 11.13 Estimators satisfying the least squares criterion for the GGLM are described as *exact weighted least squares* estimators.

The BLUE and likelihood results involve second-moment properties. Hence the forms of the result differ between the GLM and the GGLM and require separate proofs.

Lemma 11.7 For a $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{D})$ with $r = \text{rank}(\mathbf{X}) \leq q$, $\mathbf{D} = \mathbf{D}'$ known, positive definite and $N \times N$, and estimable scalar $\theta = \mathbf{c}'\boldsymbol{\beta}$, the unique BLUE of θ is $\hat{\theta} = \mathbf{c}'\tilde{\boldsymbol{\beta}}$ in which $\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}^{-1}\mathbf{y}$.

Proof. The approach is to construct and discover $\hat{\theta}$ which must be a linear estimator, namely $\hat{\theta} = \mathbf{a}'\mathbf{y}$ for some $\mathbf{a} \in \mathbb{R}^N$. The estimator $\hat{\theta}$ must be unbiased, $E(\hat{\theta}) = \theta$. If $E(\mathbf{a}'\mathbf{y}) = \mathbf{c}'\boldsymbol{\beta} \forall \boldsymbol{\beta}$, then $\mathbf{a}'\mathbf{X}\boldsymbol{\beta} = \mathbf{c}'\boldsymbol{\beta} \forall \boldsymbol{\beta}$, which implies $\mathbf{a}'\mathbf{X} = \mathbf{c}'$ and gives q linear restrictions on \mathbf{a} . That θ is estimable implies $\mathbf{c}' = \mathbf{A}\mathbf{X}$ for some \mathbf{A} ($1 \times N$). Hence the equations $\mathbf{a}'\mathbf{X} = \mathbf{c}'$ are consistent and there exists a value of \mathbf{a}' which makes $\hat{\theta}$ unbiased. Also, $\hat{\theta}$ must have minimum variance, so $\mathcal{V}(\hat{\theta}) \equiv \sigma^2\mathbf{a}'\mathbf{D}\mathbf{a}$ must be minimized subject to consistent restrictions $\mathbf{a}'\mathbf{X} = \mathbf{c}'$. The solution is found with the Lagrange function $s^2 = \sigma^2\mathbf{a}'\mathbf{D}\mathbf{a} - 2(\mathbf{a}'\mathbf{X} - \mathbf{c}')\boldsymbol{\lambda}$ for $\boldsymbol{\lambda} \in \mathbb{R}^q$. The stationary point (value of \mathbf{a}) is specified by the requirements $\partial s^2 / \partial \boldsymbol{\lambda} = \mathbf{0}$ and $\partial s^2 / \partial \mathbf{a} = \mathbf{0}$, which imply

$$\mathbf{a}'\mathbf{X} = \mathbf{c}' \tag{11.97}$$

$$\mathbf{X}\boldsymbol{\lambda} = \sigma^2 \mathbf{D}\mathbf{a}. \tag{11.98}$$

In turn, substituting $\mathbf{a} = \sigma^{-2}\mathbf{D}^{-1}\mathbf{X}\boldsymbol{\lambda}$ into the first system gives

$$\boldsymbol{\lambda}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})\sigma^{-2} = \mathbf{c}', \tag{11.99}$$

which is a consistent system. Therefore solutions are of the form

$$\boldsymbol{\lambda}' = \sigma^{-2}\mathbf{c}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-}. \tag{11.100}$$

Substituting the $\boldsymbol{\lambda}$ solution back into the original form and solving for \mathbf{a} give

$$\mathbf{a}' = \mathbf{c}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-}\mathbf{X}'\mathbf{D}^{-1}. \tag{11.101}$$

Hence $\boldsymbol{\lambda}$ and \mathbf{a} satisfy the original conditions, and $(\sigma^{-2})^{-1}$ cancels σ^{-2} .

It must now be determined that the stationary point corresponds to a minimum rather than a saddle point, which may be done by comparing $\widehat{\boldsymbol{\theta}} = \boldsymbol{\alpha}'\mathbf{y}$ with possible competitors. If $\widetilde{\boldsymbol{\theta}} = \boldsymbol{\alpha}'\mathbf{y}$ is also a LUE, then $\boldsymbol{\alpha}$ satisfies the restrictions $\boldsymbol{\alpha}'\mathbf{X} = \mathbf{c}'$. The variance of $\widetilde{\boldsymbol{\theta}}$ can be expressed as

$$\begin{aligned} \nu(\widetilde{\boldsymbol{\theta}}) &= \nu\left[(\widetilde{\boldsymbol{\theta}} - \widehat{\boldsymbol{\theta}}) + \widehat{\boldsymbol{\theta}}\right] \\ &= \nu(\widetilde{\boldsymbol{\theta}} - \widehat{\boldsymbol{\theta}}) + \nu(\widehat{\boldsymbol{\theta}}) + 2\nu\left[\widehat{\boldsymbol{\theta}}, (\widetilde{\boldsymbol{\theta}} - \widehat{\boldsymbol{\theta}})\right], \end{aligned} \tag{11.102}$$

in which

$$\begin{aligned} \nu\left[\widehat{\boldsymbol{\theta}}, (\widetilde{\boldsymbol{\theta}} - \widehat{\boldsymbol{\theta}})\right] &= \nu[\mathbf{a}'\mathbf{y}, (\boldsymbol{\alpha}' - \mathbf{a}')\mathbf{y}] \\ &= \sigma^2\mathbf{a}'\mathbf{D}(\boldsymbol{\alpha} - \mathbf{a}) \\ &= [\mathbf{c}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-}\mathbf{X}'\mathbf{D}^{-1}]\mathbf{D}(\boldsymbol{\alpha} - \mathbf{a})\sigma^2 \\ &= \mathbf{c}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-}(\mathbf{X}'\boldsymbol{\alpha} - \mathbf{X}'\mathbf{a})\sigma^2 \\ &= \mathbf{c}'(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-}(\mathbf{c} - \mathbf{c})\sigma^2 = 0. \end{aligned} \tag{11.103}$$

Hence $\widehat{\boldsymbol{\theta}}$ has minimum variance because $\nu(\widetilde{\boldsymbol{\theta}}) - \nu(\widehat{\boldsymbol{\theta}}) = \nu(\widetilde{\boldsymbol{\theta}} - \widehat{\boldsymbol{\theta}}) \geq 0$. □

Theorem 11.23 Model 1, $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{D})$ with Gaussian errors, has $r = \text{rank}(\mathbf{X}) \leq q$, $\mathbf{D} = \mathbf{D}'$ known and positive definite, and $\mathbf{D}^{-1} = \mathbf{L}\mathbf{L}'$ for nonsingular \mathbf{L} ($N \times N$). Model 2, $\text{GLM}_{N,q}[\text{row}_i(\mathbf{L}'\mathbf{y}); \text{row}_i(\mathbf{L}'\mathbf{X})\boldsymbol{\beta}, \sigma^2]$, has $\text{rank}(\mathbf{L}'\mathbf{X}) = r$, is linearly equivalent, and satisfies the least squares and Gaussian assumptions of the GLM.

Proof. $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{D})$ implies $\mathbf{L}'\mathbf{y} = \mathbf{L}'\mathbf{X}\boldsymbol{\beta} + \mathbf{L}'\boldsymbol{\epsilon} \sim \mathcal{N}_N(\mathbf{L}'\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I})$. □

Theorem 11.24 $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{D})$ with Gaussian errors, $r = \text{rank}(\mathbf{X}) \leq q$, \mathbf{D} known, positive definite, and symmetric, has estimable secondary parameter

$\theta = C\beta$. Joint MLEs of $\tau' = [\beta' \sigma^2]$ and θ are

$$\tilde{\beta} = (X'D^{-1}X)^{-1}X'D^{-1}y \tag{11.104}$$

$$\tilde{\sigma}^2 = (y - X\tilde{\beta})'D^{-1}(y - X\tilde{\beta})/N \tag{11.105}$$

and

$$\hat{\theta} = C\tilde{\beta}, \tag{11.106}$$

with $\hat{\theta}$ and $\hat{\sigma}^2$ invariant to $\tilde{\beta}$. Unbiased $\hat{\sigma}^2 = \tilde{\sigma}^2 N / (N - r)$ is usually preferred.

Proof. Left as an exercise.

Theorem 11.25 For $GGLM_{N,q}(y; X\beta, \sigma^2 D)$, if $D = D'$ is known and positive definite while $\theta = C\beta$ is estimable, the following hold.

- (a) If $\text{rank}(X) = q$, then the unique BLUE of β is $\hat{\beta} = (X'D^{-1}X)^{-1}X'D^{-1}y$.
- (b) If $\text{rank}(X) < q$, then the class of LUEs of β is empty and β is not estimable.
- (c) The unique BLUE of θ is $\hat{\theta} = C\tilde{\beta}$ in which $\tilde{\beta} = (X'D^{-1}X)^{-}X'D^{-1}y$.

Proof. Left as an exercise.

Theorem 11.26 Weighted least squares (WLS) estimators, also known as generalized least squares (GLS) estimators, may be found for $GGLM_{N,q}(y; X\beta, \sigma^2 D)$ with $r = \text{rank}(X) \leq q$ and known positive-definite, matrix D such that $D^{-1} = LL'$.

(a) Any solution $\tilde{\beta}$ of the WLS equations $(X'D^{-1}X)\tilde{\beta} = X'D^{-1}y$ is a WLS estimator for β . The solutions are of the form

$$\tilde{\beta} = (X'D^{-1}X)^{-}X'D^{-1}y \tag{11.107}$$

with

$$\begin{aligned} SSE(\tilde{\beta}) &= y'[I - L'X(X'D^{-1}X)^{-}X'L]y \\ &= y'L[I - L'X(X'D^{-1}X)^{-}X'L]L'y, \end{aligned} \tag{11.108}$$

which is invariant to the choice of generalized inverse.

(b) If $\text{rank}(X) = q$, then $(X'D^{-1}X)^{-} = (X'D^{-1}X)^{-1}$ and $\tilde{\beta} = \hat{\beta} = (X'D^{-1}X)^{-1}X'D^{-1}y$ is unique and unbiased.

Proof. Left as an exercise.

Theorem 11.27 Ordinary least squares (OLS) and WLS estimators may coincide. If

$$\hat{\beta}_D = (X'D^{-1}X)^{-1}X'D^{-1}y, \tag{11.109}$$

$$\hat{\beta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}, \tag{11.110}$$

$$\hat{\sigma}_D^2 = (\mathbf{y} - \mathbf{X}\hat{\beta}_D)' \mathbf{D}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta}_D)/N, \tag{11.111}$$

and

$$\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\hat{\beta})' \mathbf{D}^{-1}(\mathbf{y} - \mathbf{X}\hat{\beta})/N, \tag{11.112}$$

then $\hat{\beta}_D = \hat{\beta} \Leftrightarrow \exists \mathbf{B}^{-1}$ such that $\mathbf{D}\mathbf{X}\mathbf{B}^{-1} = \mathbf{X}$. If so, then $\hat{\sigma}^2 = \hat{\sigma}_D^2$.

Proof. (\Rightarrow) The fact that $\hat{\beta}_D = \hat{\beta} \ (\forall \mathbf{y})$ implies

$$\begin{aligned} (\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}^{-1} &= (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ \mathbf{X}' &= [(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{X}'\mathbf{D} \\ \mathbf{B}^{-1} &= [(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})(\mathbf{X}'\mathbf{X})^{-1}]'. \end{aligned} \tag{11.113}$$

(\Leftarrow) If such \mathbf{B} exists, then together $\mathbf{D}^{-1}\mathbf{X} = \mathbf{X}\mathbf{B}^{-1}$ and $(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1} = \mathbf{B}(\mathbf{X}'\mathbf{X})^{-1}$ give $\mathbf{D}^{-1}\mathbf{X}(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1} = \mathbf{X}\mathbf{B}^{-1}\mathbf{B}(\mathbf{X}'\mathbf{X})^{-1}$. In turn, $(\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}^{-1} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, or $\hat{\beta}_D = \hat{\beta} \ \forall \mathbf{y}$. \square

The special case of compound symmetric covariance illustrates the simplification that can occur with known covariance in a GGLM. The importance of the next lemma lies in the “univariate” approach to repeated measures and is discussed in detail in the context of hypothesis testing.

Lemma 11.8 A $\text{GGLM}_{N,q}(\mathbf{y}; \mathbf{X}\beta, \mathbf{D})$ may have \mathbf{D} which is compound symmetric (Lemma 1.33 summarizes properties),

$$\begin{aligned} \mathbf{D} &= \sigma^2[\rho\mathbf{1}_N\mathbf{1}'_N + (1 - \rho)\mathbf{I}_N] \\ &= \sigma^2 \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix} \\ &= \mathbf{V}\text{Dg}(\lambda)\mathbf{V}', \end{aligned} \tag{11.114}$$

with unknown $\sigma^2 > 0$ and $-1/(p-1) < \rho < 1$. Such models allow an exact transformation to a model with uncorrelated but heteroscedastic observations, $\text{GGLM}_{N,q}[\mathbf{V}'\mathbf{y}; \mathbf{V}'\mathbf{X}\beta, \text{Dg}(\lambda)]$.

Proof. The restrictions on the unknown parameters are necessary and sufficient for \mathbf{D} to be positive definite. The eigenvalues are $\lambda_1 = \sigma^2(1 - \rho + p\rho)$ and, $\forall j \neq 1, \lambda_j = \lambda_2 = \sigma^2(1 - \rho)$. By Lemma 1.33, the eigenvectors $\mathbf{V} = [\mathbf{v}_0 \ \mathbf{v}_1 \ \cdots \ \mathbf{v}_{N-1}]$ may be taken to be $\mathbf{v}_0 = p^{-1/2}\mathbf{1}_N$ and $[\mathbf{v}_2 \ \cdots \ \mathbf{v}_{N-1}] = \mathbf{V}_T$, the set of $N - 1$ normalized trends for N measurements. Hence \mathbf{V} may be specified

without knowing σ^2 or ρ . The result gives an exact transformation of $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$ to a model with uncorrelated but heteroscedastic observations, $\mathbf{V}'\mathbf{y} = \mathbf{V}'\mathbf{X}\boldsymbol{\beta} + \mathbf{V}'\mathbf{e}$. The new model has $\mathbf{V}'\mathbf{e} \sim \mathcal{N}_N[\mathbf{0}, \text{Dg}(\boldsymbol{\lambda})]$, corresponding to $\text{GGLM}_{N,q}[\mathbf{V}'\mathbf{y}; \mathbf{V}'\mathbf{X}\boldsymbol{\beta}, \text{Dg}(\boldsymbol{\lambda})]$. \square

EXERCISES

Prove Theorem 11.9. For a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors, estimators

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (11.49)$$

$$\widehat{\sigma}^2 = \frac{1}{N-q}\mathbf{y}'\left[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\right]\mathbf{y} \quad (11.49)$$

are

11.1 consistent,

11.2 efficient,

11.3 unbiased,

11.4 complete,

11.5 sufficient, and

11.6 UMVUE.

11.7 Furthermore $\widehat{\boldsymbol{\beta}}$ and $\widehat{\sigma}^2$ are mutually independent,

11.8 $\widehat{\boldsymbol{\beta}} \sim \mathcal{N}_q[\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]$ and

11.9 $\widehat{\sigma}^2(N-q)/\sigma^2 \sim \chi^2(N-q)$.

Hints. You may prove the results in any order you wish.

Among asymptotically unbiased estimators, one with minimum variance in large samples is called an efficient estimator or simply efficient.

You may cite results in Chapters 1–11 for the exercises. When you cite a result from Chapters 1–11, clearly indicate the number of the theorem, corollary, or lemma that you are using.

For your own enlightenment (but not for the exercises), you may wish to consider proving the results without using any results from Chapter 11.

CHAPTER 12

Estimation for Multivariate Linear Models

12.1 ALTERNATE FORMULATIONS OF THE MODEL

The univariate general linear model concerns an $N \times 1$ vector of responses \mathbf{y} with all observations independent. The multivariate general linear model allows generalizing the response vector to an $N \times p$ matrix \mathbf{Y} with some observations independent and some dependent (correlated). Only particular patterns of correlation meet the assumptions. *For data with correlated observations, distinguishing between the observational unit and the independent sampling unit often provides the key first step in choosing a valid model.*

Depending on the task at hand, it may be more convenient to express $N \times p$ \mathbf{Y} in terms of its N rows $\{\mathbf{Y}_i\}$ or its p columns $\{\mathbf{y}_j\}$. Each row, sometimes indicated $\mathbf{Y}_i = \text{row}_i(\mathbf{Y})$, is $1 \times p$, while each column, say $\mathbf{y}_j = \text{col}_j(\mathbf{Y})$, is $N \times 1$. Each row corresponds to a particular independent sampling unit, with elements within a row being observations for the sampling unit. In a study with (unrelated) human participants, each row of \mathbf{Y} contains the data for one person, while different columns contain different response variables, which might be repeated measures of a single variable. It helps to remember that

$$\mathbf{Y} = [\mathbf{y}_1 \ \mathbf{y}_2 \ \cdots \ \mathbf{y}_p] = \begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \\ \vdots \\ \mathbf{Y}_N \end{bmatrix}. \quad (12.1)$$

The first form decomposes \mathbf{Y} into p variables and the second into N independent sampling units.

The present chapter centers on deriving estimators for \mathbf{B} and $\mathbf{\Sigma}$ which satisfy a variety of optimal properties. The great majority of such properties do not depend on a particular choice of distribution function for the responses and are exact, even in small samples. In contrast, the desire to test hypotheses leads to describing distributions of test statistics. Finding exact test distributions for small samples usually requires explicit and particular specification of the distribution of the data. The most common choice involves assuming each row of errors independently follows a multivariate Gaussian distribution.

Describing the model in vector form helps explain the many relationships among the elements. The vector form helps in understanding the connections and differences among the multivariate and univariate GLM, mixed models, and more general forms. It also provides a convenient basis for many proofs of multivariate GLM properties. In parallel to expressions for \mathbf{Y} in terms of rows or columns,

$$\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_q] = \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix}, \quad (12.2)$$

$$\mathbf{B} = [\boldsymbol{\beta}_1 \ \boldsymbol{\beta}_2 \ \cdots \ \boldsymbol{\beta}_p] = \begin{bmatrix} \mathbf{B}_1 \\ \mathbf{B}_2 \\ \vdots \\ \mathbf{B}_q \end{bmatrix}, \quad (12.3)$$

and

$$\mathbf{E} = [\mathbf{e}_1 \ \mathbf{e}_2 \ \cdots \ \mathbf{e}_p] = \begin{bmatrix} \mathbf{E}_1 \\ \mathbf{E}_2 \\ \vdots \\ \mathbf{E}_N \end{bmatrix}. \quad (12.4)$$

As discussed in Chapter 3, a multivariate linear model is often written as $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$. The notation just described allows stacking the data *by variable* (column of \mathbf{Y} , \mathbf{B} , and \mathbf{E}) to give a single column of responses and errors:

$$\begin{aligned} \text{vec}(\mathbf{Y}) &= \text{vec}(\mathbf{X}\mathbf{B}) + \text{vec}(\mathbf{E}) \\ \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_p \end{bmatrix} &= \begin{bmatrix} \mathbf{X}\boldsymbol{\beta}_1 \\ \mathbf{X}\boldsymbol{\beta}_2 \\ \vdots \\ \mathbf{X}\boldsymbol{\beta}_p \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_p \end{bmatrix} \\ \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_p \end{bmatrix} &= (\mathbf{I}_p \otimes \mathbf{X})\text{vec}(\mathbf{B}) + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_p \end{bmatrix} \\ \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_p \end{bmatrix} &= (\mathbf{I}_p \otimes \mathbf{X}) \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \vdots \\ \boldsymbol{\beta}_p \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_p \end{bmatrix} \\ \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_p \end{bmatrix} &= \begin{bmatrix} \mathbf{X} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{X} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{X} \end{bmatrix} \begin{bmatrix} \boldsymbol{\beta}_1 \\ \boldsymbol{\beta}_2 \\ \vdots \\ \boldsymbol{\beta}_p \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_p \end{bmatrix}. \end{aligned} \quad (12.5)$$

Stacking the data by variable allows concluding that the submodel for response variable j is a (valid) GLM corresponding to $\mathbf{y}_j = \mathbf{X}\boldsymbol{\beta}_j + \mathbf{e}_j$, with a common

design matrix and distinct parameters $\text{GLM}_{N,q}(y_{ij}; \mathbf{X}_i\boldsymbol{\beta}_j, \sigma_j^2)$. Elements of \mathbf{X} and $\boldsymbol{\beta}_j$, combine to describe differences *between* sampling units in response variable j .

Alternatively, stacking the data by independent sampling unit (row of \mathbf{Y} , \mathbf{X} , and \mathbf{E}) focuses on differences *within* sampling unit, between correlated response variables, as captured by \mathbf{B}_k . The approach gives

$$\begin{aligned} \text{vec}(\mathbf{Y}') &= \text{vec}[(\mathbf{X}\mathbf{B})'] + \text{vec}(\mathbf{E}') \\ \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_N \end{bmatrix} &= \begin{bmatrix} (\mathbf{X}_1\mathbf{B})' \\ (\mathbf{X}_2\mathbf{B})' \\ \vdots \\ (\mathbf{X}_N\mathbf{B})' \end{bmatrix} + \begin{bmatrix} \mathbf{E}'_1 \\ \mathbf{E}'_2 \\ \vdots \\ \mathbf{E}'_N \end{bmatrix} \\ \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_N \end{bmatrix} &= (\mathbf{X} \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}') + \begin{bmatrix} \mathbf{E}'_1 \\ \mathbf{E}'_2 \\ \vdots \\ \mathbf{E}'_N \end{bmatrix} \\ \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_N \end{bmatrix} &= (\mathbf{X} \otimes \mathbf{I}_p) \begin{bmatrix} \mathbf{B}'_1 \\ \mathbf{B}'_2 \\ \vdots \\ \mathbf{B}'_q \end{bmatrix} + \begin{bmatrix} \mathbf{E}'_1 \\ \mathbf{E}'_2 \\ \vdots \\ \mathbf{E}'_N \end{bmatrix} \\ \begin{bmatrix} \mathbf{Y}'_1 \\ \mathbf{Y}'_2 \\ \vdots \\ \mathbf{Y}'_N \end{bmatrix} &= \begin{bmatrix} x_{11}\mathbf{I}_p & x_{12}\mathbf{I}_p & \cdots & x_{1q}\mathbf{I}_p \\ x_{21}\mathbf{I}_p & x_{22}\mathbf{I}_p & \cdots & x_{2q}\mathbf{I}_p \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1}\mathbf{I}_p & x_{N2}\mathbf{I}_p & \cdots & x_{Nq}\mathbf{I}_p \end{bmatrix} \begin{bmatrix} \mathbf{B}'_1 \\ \mathbf{B}'_2 \\ \vdots \\ \mathbf{B}'_q \end{bmatrix} + \begin{bmatrix} \mathbf{E}'_1 \\ \mathbf{E}'_2 \\ \vdots \\ \mathbf{E}'_N \end{bmatrix}. \end{aligned} \tag{12.6}$$

The corresponding equation for sampling unit i , namely

$$\mathbf{Y}'_i = (\mathbf{X}_i \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}') + \mathbf{E}'_i, \tag{12.7}$$

does not describe a valid univariate GLM because $\mathcal{V}(\mathbf{Y}'_i) = \mathcal{V}(\mathbf{E}'_i) = \boldsymbol{\Sigma}$. However, it does implicitly define $\text{GGLM}_{p,q}[\mathbf{Y}'_i; (\mathbf{X}_i \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}'), \boldsymbol{\Sigma}]$ for independent sampling unit i and $\text{GGLM}_{N,p,q}[\text{vec}(\mathbf{Y}'); (\mathbf{X} \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}'), \mathbf{I}_p \otimes \boldsymbol{\Sigma}]$ for all of the data. Later discussions of mixed models will demonstrate that the equation also defines a corresponding particular mixed model.

As a special case of a mixed model, the defining characteristics of a multivariate GLM are “Kronecker design” and “Kronecker covariance.” The first requires a common design matrix for all response variables (columns of \mathbf{Y} , which may be repeated measures), and the second requires a common covariance matrix for all independent sampling units (rows of \mathbf{Y} , which may be persons). Many current approaches to mixed models were developed largely to allow relaxing the restrictive assumptions of homogenous design and homogenous covariance inherent in the multivariate GLM.

Definition 12.1 (a) The matrices B ($q \times p$) and Σ ($p \times p$) are the *primary parameters* of a multivariate $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with or without Gaussian errors.

(b) Any (finite) known constant C ($a \times q$), (finite) known constant U ($p \times b$), and (finite) known constant Θ_0 ($a \times b$) define *secondary parameter* $\Theta = CBU + \Theta_0$ ($a \times b$).

Definition 12.2 (a) *Estimators* of primary parameters with good properties, especially unbiasedness, are indicated by \hat{B} and $\hat{\Sigma}$, while ones with distinct and possibly fewer desirable properties are indicated by \tilde{B} and $\tilde{\Sigma}$.

(b) Estimators of secondary parameters take the form $\hat{\Theta} = C\hat{B}U + \Theta_0$, or $\tilde{\Theta} = C\tilde{B}U + \Theta_0$, or $\hat{\Theta} = C\tilde{B}U + \Theta_0$. The third form is used only when $\hat{\Theta}$ possesses a desirable property not shared by \tilde{B} , such as unbiasedness.

(c) Covariance matrices such as $\mathcal{V}[\text{vec}(\hat{\Theta})]$ are also secondary parameters.

Definition 12.3 (a) With T_1 ($m_1 \times p$) and T_2 ($m_2 \times p$) known constant matrices, $GLM_{N,p,q_1}(Y_i; X_{i1} B_1 T_1, \Sigma)$ and $GLM_{N,p,q_2}(Y_i; X_{i2} B_2 T_2, \Sigma)$ are *linearly equivalent between subjects* whenever the *columns* of X_1 and X_2 span the same subspace of \mathbb{R}^q .

(b) The two models are *linearly equivalent within subject* whenever the *rows* of T_1 and T_2 span the same subspace of \mathbb{R}^p .

(c) If both conditions **(a)** and **(b)** hold, then the two models are simply *linearly equivalent*.

Linear equivalence describes the expected values, the means, of $\{y_{ij}\}$, because $E(Y) = X_1 B_1 T_1$ and $E(Y) = X_2 B_2 T_2$. While X_1 and X_2 provide between-subject information, T_1 and T_2 provide within-subject information. Most often, one or both of T_1 and T_2 are simply I_p (invisible). With linearly equivalent between-subject models, (1) for any B_1 there exists B_2 such that $X_1 B_1 = X_2 B_2$ and (2) for any B_2 there exists B_1 such that $X_1 B_1 = X_2 B_2$. With linearly equivalent within-subject models, (1) for any B_1 there exists B_2 such that $B_1 T_1 = B_2 T_2$ and (2) for any B_2 there exists B_1 such that $B_1 T_1 = B_2 T_2$.

Obviously, if $T_1 = T_2$, then no attention must be paid to the question of linear equivalence within subjects. Often we have $T_1 = T_2 = I_p$. In such cases, verifying linear equivalence of two multivariate models ($p > 1$) reduces to considering only between-subject properties. The multivariate setting requires the univariate requirement to apply simultaneously to p parameter vector pairs. Here

$$Y = [y_1 \ y_2 \ \cdots \ y_p] \tag{12.8}$$

$$B_1 = [\beta_{1,1} \ \beta_{2,1} \ \cdots \ \beta_{p,1}] \tag{12.9}$$

$$B_2 = [\beta_{1,2} \ \beta_{2,2} \ \cdots \ \beta_{p,2}]. \tag{12.10}$$

Linear equivalence of $Y = X_1 B_1 + E_1$ and $Y = X_2 B_2 + E_2$ corresponds to

linear equivalence of all p univariate model pairs; $\mathbf{y}_j = \mathbf{X}_1\boldsymbol{\beta}_{j,1} + \mathbf{e}_{j,1}$ must be linearly equivalent to $\mathbf{y}_j = \mathbf{X}_2\boldsymbol{\beta}_{j,2} + \mathbf{e}_{j,2}$.

If $\mathbf{X}_1 = \mathbf{X}_2$, then no question arises about linear equivalence between subjects. Similarly, if $\mathbf{T}_1 = \mathbf{T}_2$, then no question arises about linear equivalence within subjects. If both conditions hold, the term “linear equivalence” can be discussed without ambiguity.

Example 12.1 A growth curve model provides one application of $\mathbf{T}_1 \neq \mathbf{I}_p$, with the columns of \mathbf{Y} containing repeated measures and \mathbf{T} defining the (within-subject) predictor values of interest as a function of time. With $p = 3$ times and no between-subject factor, a simple example is

$$\begin{aligned} \mathbf{Y} &= \mathbf{XBT} + \mathbf{E} \\ &= \mathbf{1}_N[\beta_0 \beta_1 \beta_2] \begin{bmatrix} 1 & 1 & 1 \\ t_1 & t_2 & t_3 \\ t_1^2 & t_2^2 & t_3^2 \end{bmatrix} + \mathbf{E}. \end{aligned} \tag{12.11}$$

The corresponding scalar form is $y_{ij} = \beta_0 + \beta_1 t_1 + \beta_2 t_1^2 + e_{ij}$. The model may be converted to one without \mathbf{T} :

$$\begin{aligned} \mathbf{YT}^{-1} &= \mathbf{XBT}^{-1} + \mathbf{ET}^{-1} \\ \mathbf{Y}_T &= \mathbf{XB}_T + \mathbf{E}_T. \end{aligned} \tag{12.12}$$

If \mathbf{T} has $m \leq p$ rows, eliminating \mathbf{T} requires multiplying by $\mathbf{T}'(\mathbf{TT}')^{-1} = \mathbf{T}^+$.

Overall, four classes of multivariate GLMs are defined by allowing either FR or LTFR designs, combined with either Gaussian or unspecified distributions for responses. For $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}|\mathbf{R}_x\mathbf{B}\mathbf{R}_y = \mathbf{A}, \boldsymbol{\Sigma})$, with or without Gaussian errors, estimation theory for primary parameter \mathbf{B} and secondary parameter $\boldsymbol{\Theta}$ is closely tied to $r = \text{rank}([\mathbf{X}' \mathbf{R}'_x]')$.

The distinction between FR and LTFR describes properties of the columns of $[\mathbf{X}' \mathbf{R}'_x]'$, which correspond directly to properties of the rows of \mathbf{B} . The form of the model statement implies the properties apply only to every element of each row and hence to each column in \mathbf{B} considered separately. Just as for the univariate models considered in the preceding chapter, every multivariate GLM constrained by $\mathbf{R}_x \neq \mathbf{I}_q$ has a linearly equivalent unconstrained model.

The presence of $\mathbf{R}_y \neq \mathbf{I}_p$ describes constraints among the columns of \mathbf{Y} , which correspond directly to properties of the columns of \mathbf{B} . Such constraints may either introduce or eliminate singularity of the error covariance matrix. Not surprisingly, every multivariate GLM constrained by $\mathbf{R}_y \neq \mathbf{I}_p$ has a linearly equivalent unconstrained model.

In the absence of \mathbf{R}_x , which gives row restrictions on \mathbf{B} , the distinction between FR and LTFR depends solely on $r = \text{rank}(\mathbf{X})$. FR models have an unbiased estimator for \mathbf{B} and for all $\boldsymbol{\Theta}$, while LTFR models never have an unbiased estimator for \mathbf{B} and have an unbiased estimator for only some $\boldsymbol{\Theta}$.

Whether or not the model is FR, unbiased estimators of Σ and related properties are always available [as long as $N > r = \text{rank}(\{\mathbf{X}' \mathbf{R}'_x\}')$]. The presence or absence of \mathbf{R}_y has no effect on bias of estimators. However, it may introduce singularity and hence cause distributions to degenerate. In turn, hypothesis tests may be undefined and require modification. Just as an equivalent FR model can be found for any LTFR model, an equivalent “nonsingular” model can be found for any “singular” model.

Theorem 12.1 Every multivariate GLM with LTFR design matrix has a linearly equivalent model with a full-rank design matrix. In particular, for $\text{GLM}_{N,p,q}\text{LTFR}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with $\text{rank}(\mathbf{X}) = r < q$, there exists $\text{GLM}_{N,p,r}\text{FR}(\mathbf{Y}_i; \mathbf{X}_{i1} \mathbf{B}_1, \Sigma)$ which is linearly equivalent with $\text{rank}(\mathbf{X}_1) = r$.

Proof. With the subscript 1 indicating the components corresponding to positive singular values, the SVD allows writing

$$\begin{aligned} \mathbf{X} &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{R}'_1 \\ \mathbf{R}'_0 \end{bmatrix} \\ &= \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \mathbf{R}'_1 \end{aligned} \tag{12.13}$$

and also allows defining

$$\mathbf{X}_1 = \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \tag{12.14}$$

$$\mathbf{B}_1 = \mathbf{R}'_1 \mathbf{B}. \tag{12.15}$$

In turn

$$\begin{aligned} \mathbf{Y} &= \mathbf{X} \mathbf{B} + \mathbf{E} \\ &= \mathbf{X}_1 \mathbf{B}_1 + \mathbf{E}, \end{aligned} \tag{12.16}$$

with full rank \mathbf{X}_1 . Although well behaved and elegant, the full-rank model is merely one particular choice among infinitely many equivalent models. \square

12.2 ESTIMABILITY IN THE MULTIVARIATE GLM

For secondary parameter $\Theta = \mathbf{C} \mathbf{B} \mathbf{U}$, the within-subject contrast matrix \mathbf{U} helps define Θ . It also plays a key role in determining (within-subject) linear equivalence as well as whether or not a valid hypothesis test exists. However, it has no effect on estimability.

Definition 12.4 A $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_{i1} \mathbf{B}, \Sigma)$ has $\Theta = \mathbf{C} \mathbf{B} \mathbf{U}$.

- (a) Primary parameter \mathbf{B} is *estimable* if and only if a $q \times N$ constant matrix \mathbf{A}_1 exists such that $\text{E}(\mathbf{A}_1 \mathbf{Y}) = \mathbf{B}$.
- (b) Secondary parameter Θ is *estimable* if and only if constant matrix \mathbf{A}_2 ($a \times N$) exists such that $\text{E}(\mathbf{A}_2 \mathbf{Y} \mathbf{U}) = \Theta$.
- (c) For known, fixed Θ_0 , $\Theta + \Theta_0$ is *estimable* if and only if Θ is *estimable*.

Part (a) of the definition is essentially redundant since it is a special case of part (b). Choosing $C = I_q$ defines B as a secondary parameter.

Theorem 12.2 In $GLM_{N,p,q}FR(Y_i; X_i B, \Sigma)$ with $\Theta = CBU$ defined by constants C and U , both B and Θ are estimable.

Proof. The proof is essentially the same as for univariate models. Only the properties of the design matrix X are involved in the issue of estimability. \square

Theorem 12.3 For $GLM_{N,p,q}LTFR(Y_i; X_i B, \Sigma)$ with $\Theta = CBU$, the following hold.

- (a) With A_1 $q \times N$ and constant, *no* linear transformation of the data, say $\widehat{B} = A_1 Y$, exists such that $E(\widehat{B}) = B$.
- (b) For $a \times b$ secondary parameter Θ , with constants A_2 $a \times N$ and U $p \times b$, a linear transformation of the data, say $\widehat{\Theta} = A_2 Y U$, *may* or *may not* exist such that $E(\widehat{\Theta}) = \Theta$.
- (c) If $X_i = \text{row}_i(X)$, $E_i = \text{row}_i(I_N)$, and $\theta_{ij} = E(E_i Y E_j') = E_i X B E_j' = X_i B E_j$, then θ_{ij} is always *estimable*.
- (d) The within-subject contrast matrix U plays no role in determining estimability of Θ , only C does. Secondary parameters $\Theta = CBU$ and $\Theta - \Theta_0$ are estimable if and only if CB is estimable.

Proof. Essentially the same as for univariate models, by considering individual columns of U . Only the properties of X are involved in the issue of estimability.

Theorem 12.4 Any primary or secondary expected-value parameter estimable in $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ is also estimable in a linearly equivalent model.

Proof. Considering estimable $\Theta = CBU$ allows treating B as a special case. Extending the known result for $p = 1$ (the univariate result) is easy due to the unimportance of U in determining estimability. Details are left as an exercise.

Estimators for the univariate GLM arise naturally as special cases of the multivariate forms. However, multivariate models require expanding the concepts and results surrounding estimability and restricted models. Furthermore, differences in development are sufficient to warrant separate treatment in some cases, especially for maximum likelihood estimation. We have omitted separate development of least squares estimators for the multivariate case. The task is straightforward and can easily be done by the interested reader. Although distribution free, the key mathematical forms coincide with or closely resemble the forms needed for deriving maximum likelihood estimators.

12.3 UNRESTRICTED LIKELIHOOD ESTIMATION

Derivation of the maximum likelihood estimators of \mathbf{B} and Σ under $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ assumptions with Gaussian variables is lengthy but straightforward. In formulating the likelihood function it helps to keep in mind the close relationship between the univariate and multivariate GLM.

Lemma 12.1 Model 1, $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with constant \mathbf{C} has estimable $a \times p$ secondary parameter $\Theta = \mathbf{C}\mathbf{B}$. The assumptions and definitions of the model can also be expressed as specified in models 2 and 3 below.

Model 2, $\text{GGLM}_{N,p,q}[\text{vec}(\mathbf{Y}); (\mathbf{I}_p \otimes \mathbf{X})\text{vec}(\mathbf{B}), \Sigma \otimes \mathbf{I}_N]$ and estimable $ap \times 1$ secondary parameter

$$\begin{aligned}\tau &= \text{vec}(\Theta) \\ &= \text{vec}(\mathbf{C}\mathbf{B}) \\ &= (\mathbf{I}_p \otimes \mathbf{C})\text{vec}(\mathbf{B}).\end{aligned}\quad (12.17)$$

Model 3, $\text{GGLM}_{N,p,q}[\text{vec}(\mathbf{Y}'); (\mathbf{X} \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}'), (\mathbf{I}_N \otimes \Sigma)]$ and estimable $ap \times 1$ secondary parameter

$$\begin{aligned}\tau_T &= \text{vec}(\Theta') \\ &= \text{vec}(\mathbf{B}'\mathbf{C}') \\ &= (\mathbf{C} \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}'),\end{aligned}\quad (12.18)$$

with τ_T and τ differing only by being a permutation of each other.

Proof. Results follow directly from the definitions of the multivariate GLM, the GGLM, the $\text{vec}()$ operator, and the Kronecker product $\mathbf{A} \otimes \mathbf{B}$. In particular, Theorem 1.5 gives $\text{vec}(\mathbf{A}\mathbf{B}\mathbf{C}) = (\mathbf{C}' \otimes \mathbf{A})\text{vec}(\mathbf{B})$. \square

Lemma 12.2 For a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and $\text{rank}(\mathbf{X}) \leq q$, the log likelihood may be written

$$\log L(\mathbf{B}, \Sigma; \mathbf{Y}_*) = -N \log |2\pi\Sigma|/2 - \text{tr}[\Sigma^{-1}(\mathbf{Y}_* - \mathbf{X}\mathbf{B})'(\mathbf{Y}_* - \mathbf{X}\mathbf{B})]/2. \quad (12.19)$$

Proof. With \mathbf{Y}_i $1 \times p$ and \mathbf{X}_i $1 \times q$, the joint density function for $\text{vec}(\mathbf{Y})$ is

$$\begin{aligned}L &= f_{\mathbf{Y}}(\mathbf{Y}_*) = \prod_{i=1}^N |2\pi\Sigma|^{-1/2} \exp[-(\mathbf{Y}_{i*} - \mathbf{X}_i\mathbf{B})\Sigma^{-1}(\mathbf{Y}_{i*} - \mathbf{X}_i\mathbf{B})'/2] \\ &= |2\pi\Sigma|^{-N/2} \exp\{-\text{tr}[(\mathbf{Y}_* - \mathbf{X}\mathbf{B})\Sigma^{-1}(\mathbf{Y}_* - \mathbf{X}\mathbf{B})']/2\}.\end{aligned}\quad (12.20)$$

In turn $\log L = -N \log |2\pi\Sigma|/2 - \text{tr}[(\mathbf{Y}_* - \mathbf{X}\mathbf{B})\Sigma^{-1}(\mathbf{Y}_* - \mathbf{X}\mathbf{B})']/2$. The final form follows from the cyclical property of the trace function. \square

Lemma 12.3 For a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and $\text{rank}(\mathbf{X}) \leq q$, if $\mathbf{y}_s = \text{vec}(\mathbf{Y})$, $\mathbf{X}_s = \mathbf{I}_p \otimes \mathbf{X}$, $\beta_s = \text{vec}(\mathbf{B})$, $\Sigma_s = (\Sigma \otimes \mathbf{I}_N)$,

$\mathbf{e}_s = \text{vec}(\mathbf{E})$, and $\mathbf{E} = (\mathbf{Y} - \mathbf{X}\mathbf{B})$, then the log likelihood is

$$\begin{aligned} \log L(\boldsymbol{\beta}_s, \boldsymbol{\Sigma}_s; \mathbf{y}_{s*}) &= -\log |2\pi\boldsymbol{\Sigma}_s|/2 - (\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)' \boldsymbol{\Sigma}_s^{-1} (\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)/2 \quad (12.21) \\ &= -\log |2\pi\boldsymbol{\Sigma}_s|/2 - \text{tr}(\boldsymbol{\Sigma}_s^{-1} \mathbf{e}_{s*} \mathbf{e}'_{s*})/2. \end{aligned}$$

Proof. The joint density function of \mathbf{y}_s can be written as

$$f_{\mathbf{y}_s}(\mathbf{y}_{s*}) = |2\pi\boldsymbol{\Sigma}_s|^{-1/2} \exp[-(\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)' \boldsymbol{\Sigma}_s^{-1} (\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)/2]. \quad (12.22)$$

Further details are left to the reader. □

The following lemmas are in numerous matrix theory books and will be used in the proof of the next theorem. Proofs are left as exercises.

Lemma 12.4 (a) If \mathbf{A} is symmetric and nonsingular, then $\partial \log |\mathbf{A}| / \partial \mathbf{A} = 2\mathbf{A}^{-1} - \text{Dg}(\{\{\mathbf{A}^{-1}\}_{jj}\})$.

(b) If \mathbf{X} is symmetric and \mathbf{a} is fixed, then $\partial(\mathbf{a}'\mathbf{X}\mathbf{a}) / \partial \mathbf{X} = 2\mathbf{a}\mathbf{a}' - \text{Dg}(\mathbf{a}\mathbf{a}')$.

Lemma 12.5 (a) If \mathbf{A} , \mathbf{B} , \mathbf{C} , and \mathbf{D} conform to the operations, then $(\mathbf{A} \otimes \mathbf{B})(\mathbf{C} \otimes \mathbf{D}) = (\mathbf{AC} \otimes \mathbf{BD})$.

(b) If \mathbf{X} and \mathbf{B} conform to the operations, then $\text{vec}(\mathbf{XB}) = (\mathbf{I} \otimes \mathbf{X})\text{vec}(\mathbf{B})$ and $\text{vec}(\mathbf{B}'\mathbf{X}') = (\mathbf{X} \otimes \mathbf{I})\text{vec}(\mathbf{B}')$.

(c) If \mathbf{E} and \mathbf{S} conform to the operations, then $[\text{vec}(\mathbf{E})]'(\mathbf{S} \otimes \mathbf{I})[\text{vec}(\mathbf{E})] = [\text{vec}(\mathbf{E}')]'(\mathbf{I} \otimes \mathbf{S}')\text{vec}(\mathbf{E}')$.

(d) $(\mathbf{A} \otimes \mathbf{B})^- = \mathbf{A}^- \otimes \mathbf{B}^-$ (which is just one of infinitely many).

Theorem 12.5 For $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors, $r = \text{rank}(\mathbf{X}) \leq q$, and $a \times b$ estimable secondary parameter $\boldsymbol{\Theta} = \mathbf{C}\mathbf{B}\mathbf{U}$, the joint supremum (for \mathbf{B} if $r < q$) or maximum (if $r = q$ for \mathbf{B} ; always for $\boldsymbol{\Sigma}$ and estimable $\boldsymbol{\Theta}$) likelihood estimators of \mathbf{B} , $\boldsymbol{\Sigma}$, and $\boldsymbol{\Theta}$ are

$$\tilde{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{Y} \quad (12.23)$$

$$\tilde{\boldsymbol{\Sigma}} = (\mathbf{Y} - \mathbf{X}\tilde{\mathbf{B}})'(\mathbf{Y} - \mathbf{X}\tilde{\mathbf{B}})/N = \mathbf{Y}'[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}']\mathbf{Y}/N \quad (12.24)$$

$$\hat{\boldsymbol{\Theta}} = \mathbf{C}\tilde{\mathbf{B}}\mathbf{U}, \quad (12.25)$$

with $\hat{\boldsymbol{\Theta}}$ invariant to $\tilde{\mathbf{B}}$ for estimable $\boldsymbol{\Theta}$. Here $\tilde{\mathbf{B}}$ is any solution of $(\mathbf{X}'\mathbf{X})\tilde{\mathbf{B}} = \mathbf{X}'\mathbf{Y}$, with infinitely many for $r < q$, and one unique solution if $r = q$ and $(\mathbf{X}'\mathbf{X})^- = (\mathbf{X}'\mathbf{X})^{-1}$. The value of $\tilde{\boldsymbol{\Sigma}}$ is always invariant to $(\mathbf{X}'\mathbf{X})^-$. It is customary to use $\hat{\boldsymbol{\Sigma}} = \tilde{\boldsymbol{\Sigma}}N/(N - r)$, which is unbiased.

Proof. To obtain the solution equations, we differentiate the log likelihood with respect to the elements of \mathbf{B} and $\boldsymbol{\Sigma}$, set the derivatives to zero, and solve for \mathbf{B} and $\boldsymbol{\Sigma}$. Next $\mathbf{y}_s = \text{vec}(\mathbf{Y})$ is $Np \times 1$, $\mathbf{X}_s = \mathbf{I}_p \otimes \mathbf{X}$ ($Np \times pq$), $\boldsymbol{\beta}_s = \text{vec}(\mathbf{B})$, $\boldsymbol{\Sigma}_s = \boldsymbol{\Sigma} \otimes \mathbf{I}_N$, $\mathbf{e}_s = \text{vec}(\mathbf{E}) = (\mathbf{y}_s - \mathbf{X}_s\boldsymbol{\beta}_s)$, and $\mathbf{E} = \mathbf{Y} - \mathbf{X}\mathbf{B}$. By Lemma 12.3

$$\log L(\mathbf{B}, \Sigma; \mathbf{Y}_*) = -\log|2\pi\Sigma_s|/2 - (\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)\Sigma_s^{-1}(\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)'/2. \quad (12.26)$$

If $t_0 = -(Np/2)\log(2\pi)$, $t_1(\Sigma) = -(1/2)\log|\Sigma_s|$, and $t_2(\boldsymbol{\beta}, \Sigma) = (\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)'\Sigma_s^{-1}(\mathbf{y}_{s*} - \mathbf{X}_s\boldsymbol{\beta}_s)$, then $\log L(\mathbf{B}, \Sigma; \mathbf{Y}_*) = t_0 + t_1(\Sigma) - t_2(\mathbf{B}, \Sigma)/2$. Also,

$$\begin{aligned} t_1(\Sigma) &= -(1/2)\log|\Sigma \otimes \mathbf{I}_N| \\ &= (1/2)\log|\Sigma \otimes \mathbf{I}_N|^{-1} \\ &= (1/2)\log\left(|\Sigma^{-1}|^N |\mathbf{I}_N|^p\right) \\ &= (N/2)\log|\Sigma^{-1}|. \end{aligned} \quad (12.27)$$

Using direct-product properties in Lemma 12.5,

$$\begin{aligned} t_2(\boldsymbol{\beta}_s, \Sigma) &= \mathbf{y}'_{s*}(\Sigma \otimes \mathbf{I}_N)^{-1}\mathbf{y}_{s*} - 2\boldsymbol{\beta}'_s(\mathbf{I}_p \otimes \mathbf{X}')(\Sigma \otimes \mathbf{I}_N)^{-1}\mathbf{y}_{s*} + \\ &\quad \boldsymbol{\beta}'_s(\mathbf{I}_p \otimes \mathbf{X}')(\Sigma \otimes \mathbf{I}_N)^{-1}(\mathbf{I}_p \otimes \mathbf{X})\boldsymbol{\beta}_s \\ &= \mathbf{y}'_{s*}(\Sigma^{-1} \otimes \mathbf{I}_N)\mathbf{y}_{s*} - 2\boldsymbol{\beta}'_s(\Sigma^{-1} \otimes \mathbf{X}')\mathbf{y}_{s*} + \boldsymbol{\beta}'_s(\Sigma^{-1} \otimes \mathbf{X}'\mathbf{X})\boldsymbol{\beta}_s. \end{aligned} \quad (12.28)$$

Differentiating $\log L$ with respect to $\boldsymbol{\beta}_s$ gives

$$\begin{aligned} \partial \log L / \partial \boldsymbol{\beta}_s &= \partial[-t_2(\boldsymbol{\beta}_s, \Sigma)/2] / \partial \boldsymbol{\beta}_s \\ &= \mathbf{0} + (\Sigma^{-1} \otimes \mathbf{X}')\mathbf{y}_{s*} - (\Sigma^{-1} \otimes \mathbf{X}'\mathbf{X})\boldsymbol{\beta}_s. \end{aligned} \quad (12.29)$$

Setting the expression to zero (and using Lemma 12.5) gives the MLE for $\boldsymbol{\beta}_s$,

$$\begin{aligned} \tilde{\boldsymbol{\beta}}_s &= (\Sigma^{-1} \otimes \mathbf{X}'\mathbf{X})^{-1}(\Sigma^{-1} \otimes \mathbf{X}')\mathbf{y}_s \\ &= [\Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1}](\Sigma^{-1} \otimes \mathbf{X}')\mathbf{y}_s \\ &= [\mathbf{I}_p \otimes (\mathbf{X}'\mathbf{X})^{-1}]\mathbf{X}'\mathbf{y}_s. \end{aligned} \quad (12.30)$$

Surprisingly, and very conveniently, Σ cancels! Equivalently, $\text{vec}(\tilde{\mathbf{B}}) = [\mathbf{I}_p \otimes (\mathbf{X}'\mathbf{X})^{-1}]\text{vec}(\mathbf{Y})$. By Lemma 12.4, $\tilde{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$. In summary, $\tilde{\mathbf{B}} = [\tilde{\boldsymbol{\beta}}_1 \cdots \tilde{\boldsymbol{\beta}}_p]$ with $\tilde{\boldsymbol{\beta}}_j = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}_j$ and

$$\tilde{\boldsymbol{\beta}}_s = \text{vec}(\tilde{\mathbf{B}}) = \begin{bmatrix} \tilde{\boldsymbol{\beta}}_1 \\ \tilde{\boldsymbol{\beta}}_2 \\ \vdots \\ \tilde{\boldsymbol{\beta}}_p \end{bmatrix}. \quad (12.31)$$

The cancellation of Σ in the derivation of $\boldsymbol{\beta}_s$ is extremely important. The “disappearance” allows estimating $\boldsymbol{\beta}_s$ (and thus \mathbf{B}) in a noniterative fashion. If Σ did not disappear from the likelihood equations, one would be required to know or estimate Σ before estimating \mathbf{B} . The strong structural assumptions allow Σ to cancel. Rows of \mathbf{Y} are independent, each row of \mathbf{Y} has the same covariance matrix, and each column of \mathbf{Y} has the same design.

Since Σ and Σ^{-1} have a one-to-one correspondence, maximization with respect to Σ^{-1} is the same as maximization with respect to Σ . It is more convenient to take derivatives of $\log L$ with respect to Σ^{-1} :

$$\frac{\partial \log L}{\partial \Sigma^{-1}} = \frac{\partial}{\partial \Sigma^{-1}} t_1(\Sigma) - \frac{1}{2} \frac{\partial}{\partial \Sigma^{-1}} t_2(\beta_s, \Sigma). \quad (12.32)$$

By Lemma 12.4

$$\frac{\partial}{\partial \Sigma^{-1}} t_1(\Sigma) = \frac{N}{2} [2\Sigma - \text{Dg}(\Sigma)]. \quad (12.33)$$

For differentiating with respect to Σ^{-1} , it is more convenient to express t_2 as a sum of quadratic forms in Σ^{-1} , which corresponds to working with the likelihood function of $\text{vec}(\mathbf{Y}')$ (i.e., \mathbf{Y} stacked by rows) rather than $\text{vec}(\mathbf{Y})$. With $\mathbf{E}_* = (\mathbf{Y}_* - \mathbf{X}\mathbf{B})$ and $\mathbf{E}_{i*} = \text{row}_i(\mathbf{E}_*)$, direct-product properties in Lemma 12.5 give

$$\begin{aligned} t_2(\beta, \Sigma) &= (\mathbf{y}_{s*} - \mathbf{X}_s \beta_s)' \Sigma_s^{-1} (\mathbf{y}_{s*} - \mathbf{X}_s \beta_s) \\ &= [\text{vec}(\mathbf{E}_*)]' (\Sigma^{-1} \otimes \mathbf{I}_N) [\text{vec}(\mathbf{E}_*)] \\ &= [\text{vec}(\mathbf{E}'_*)]' (\mathbf{I}_N \otimes \Sigma^{-1}) [\text{vec}(\mathbf{E}'_*)] \\ &= \sum_{i=1}^N \mathbf{E}_{i*} \Sigma^{-1} \mathbf{E}'_{i*}. \end{aligned} \quad (12.34)$$

Using Lemma 12.4 we have

$$\begin{aligned} \frac{\partial}{\partial \Sigma^{-1}} t_2(\beta, \Sigma) &= \sum_{i=1}^N \frac{\partial}{\partial \Sigma^{-1}} \mathbf{E}_{i*} \Sigma^{-1} \mathbf{E}'_{i*} \\ &= \sum_{i=1}^N [2\mathbf{E}'_{i*} \mathbf{E}_{i*} - \text{Dg}(\{\langle \mathbf{E}'_{i*} \mathbf{E}_{i*} \rangle_{jj}\})] \\ &= 2\mathbf{E}'_* \mathbf{E}_* - \text{Dg}(\{\langle \mathbf{E}'_* \mathbf{E}_* \rangle_{jj}\}) \\ &= 2(\mathbf{Y}_* - \mathbf{X}\mathbf{B})' (\mathbf{Y}_* - \mathbf{X}\mathbf{B}) - \text{Dg}(\{\langle (\mathbf{Y}_* - \mathbf{X}\mathbf{B})' (\mathbf{Y}_* - \mathbf{X}\mathbf{B}) \rangle_{jj}\}). \end{aligned} \quad (12.35)$$

The MLEs for \mathbf{B} and Σ are produced by combining results and setting the derivatives equal to zero. Considering $\partial \log L(\mathbf{B}, \Sigma, \mathbf{Y}_{s*}) / \partial \Sigma^{-1} = \mathbf{0}$ implies

$$[2\tilde{\Sigma} - \text{Dg}(\tilde{\Sigma})]N/2 = [2\hat{\mathbf{E}}' \hat{\mathbf{E}}/N - \text{Dg}(\{\langle \mathbf{E}' \mathbf{E} \rangle_{jj}\})/N]N/2, \quad (12.36)$$

in which $\hat{\mathbf{E}} = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})$. Hence $\tilde{\Sigma} = \hat{\mathbf{E}}' \hat{\mathbf{E}}/N = (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})' (\mathbf{Y} - \mathbf{X}\hat{\mathbf{B}})/N$.

Differentiating the log likelihood function and setting the results to zeros gave $\hat{\mathbf{B}}$ and $\tilde{\Sigma}$ as solutions of the likelihood equations. The uniqueness of the solutions implies that if the likelihood function has an extremum (minimum, maximum, or saddle point), then our estimators provide the coordinates of the extremum.

The likelihood function, being a density, is nonnegative for all values of the variable (\mathbf{Y}_*) and the parameters. For a fixed value of \mathbf{Y}_* (say \mathbf{Y}_0), one can make the likelihood arbitrarily close to zero by choosing $\Sigma = \mathbf{I}_p$ and β_s such that $[\text{vec}(\mathbf{Y}_0) - \mathbf{X}_s \beta_s] \rightarrow \pm \infty$. Thus, the likelihood function has an infimum at zero but no minimum. The only possibilities are that $\hat{\mathbf{B}}$ and $\tilde{\Sigma}$ locate a maximum or a saddle point.

To demonstrate the extremum is a maximum for $\text{rank}(\mathbf{X}) = r = q$ and supremum for $r < q$, one must prove

$$\mathbf{H} = \begin{bmatrix} \partial^2 \log L / (\partial \boldsymbol{\beta}_s \partial \boldsymbol{\beta}'_s) & \partial^2 \log L / (\partial \boldsymbol{\beta}_s \partial \boldsymbol{\Sigma}^{-1}) \\ \partial^2 \log L / (\partial \boldsymbol{\Sigma}^{-1} \partial \boldsymbol{\beta}'_s) & \partial^2 \log L / (\partial \boldsymbol{\Sigma}^{-1} \partial \boldsymbol{\Sigma}^{-1}) \end{bmatrix} \quad (12.37)$$

is negative definite at $\widehat{\mathbf{B}}$ and $\widehat{\boldsymbol{\Sigma}}$. Details will be left to the reader. \square

Theorem 12.6 For $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with $r = \text{rank}(\mathbf{X}) \leq q$ and $\boldsymbol{\beta}_s = \text{vec}(\mathbf{B})$ ($pq \times 1$) the following results hold.

(a) If $r < q$, then no BLUE exists and \mathbf{B} is not estimable.

(b) If $r = q$, then $\text{BLUE}(\mathbf{B}) = \widehat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$.

(c) Equivalently $\text{BLUE}(\boldsymbol{\beta}_s) = \widehat{\boldsymbol{\beta}}_s = \text{vec}(\widehat{\mathbf{B}}) = [\mathbf{I} \otimes (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] \text{vec}(\mathbf{Y})$. (12.38)

Proof. By Lemma 12.1, the assumptions of $\text{GLM}_{N,p,q}\text{FR}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ may be expressed as $\text{GGLM}_{Np,q}\text{FR}(\mathbf{y}_s; \mathbf{X}_s \mathbf{B}_s, \boldsymbol{\Sigma}_s)$, in which \mathbf{y}_s is $Np \times 1$, $\boldsymbol{\Sigma}$ is $p \times p$ and \mathbf{X} is $N \times q$, $\mathbf{y}_s = \text{vec}(\mathbf{Y})$, $\mathbf{X}_s = (\mathbf{I}_p \otimes \mathbf{X})$, $\boldsymbol{\beta}_s = \text{vec}(\mathbf{B})$, and $\boldsymbol{\Sigma}_s = (\boldsymbol{\Sigma} \otimes \mathbf{I}_N)$. By Theorem 11.25 on weighted estimation, the BLUE of $\boldsymbol{\beta}_s$ is

$$\begin{aligned} \widehat{\boldsymbol{\beta}}_s &= (\mathbf{X}'_s \boldsymbol{\Sigma}_s^{-1} \mathbf{X}_s)^{-1} \mathbf{X}'_s \boldsymbol{\Sigma}_s^{-1} \mathbf{y}_s \\ &= [(\mathbf{I}_p \otimes \mathbf{X}')' (\boldsymbol{\Sigma} \otimes \mathbf{I}_N)^{-1} (\mathbf{I}_p \otimes \mathbf{X})]^{-1} (\mathbf{I}_p \otimes \mathbf{X}') (\boldsymbol{\Sigma} \otimes \mathbf{I}_N)^{-1} \mathbf{y}_s \\ &= [(\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}') (\mathbf{I}_p \otimes \mathbf{X})]^{-1} (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}') \mathbf{y}_s \\ &= (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}'\mathbf{X})^{-1} (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}') \mathbf{y}_s \\ &= [\boldsymbol{\Sigma} \otimes (\mathbf{X}'\mathbf{X})^{-1}] (\boldsymbol{\Sigma}^{-1} \otimes \mathbf{X}') \mathbf{y}_s \\ &= [\mathbf{I} \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'] \mathbf{y}_s \\ &= \begin{bmatrix} (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{y}_1 \\ \vdots \\ (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}' \mathbf{y}_p \end{bmatrix}. \end{aligned} \quad (12.39)$$

Hence $\widehat{\mathbf{B}}$ must satisfy $\widehat{\boldsymbol{\beta}}_s = \text{vec}(\widehat{\mathbf{B}}) = [\widehat{\boldsymbol{\beta}}'_1 \widehat{\boldsymbol{\beta}}'_2 \cdots \widehat{\boldsymbol{\beta}}'_p]'$. \square

Here $\boldsymbol{\Sigma}$ “drops out” in the derivation of $\widehat{\mathbf{B}}$. Thus, unlike the situation in weighted least squares, one need not know the value of $\boldsymbol{\Sigma}$ in order to determine the BLUE in the multivariate GLM.

Corollary 12.6 The covariance matrix of $\text{BLUE}(\mathbf{B})$ is $\mathcal{V}(\widehat{\boldsymbol{\beta}}_s) = \boldsymbol{\Sigma} \otimes (\mathbf{X}'\mathbf{X})^{-1}$.

Proof.

$$\begin{aligned}
 \mathcal{V}(\hat{\beta}_s) &= [I_p \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'] \mathcal{V}(\mathbf{y}_s) [I_p \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}']' \\
 &= [I_p \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'] (\Sigma \otimes I_N) [I_p \otimes \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \\
 &= [\Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'] [I_p \otimes \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}] \\
 &= \Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1} \\
 &= \Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1}.
 \end{aligned}
 \tag{12.40}$$

□

12.4 ESTIMATION OF SECONDARY PARAMETERS

The commonly described format of secondary parameters in the multivariate GLM is the $a \times b$ form $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$, with \mathbf{C} and \mathbf{U} usually full-rank matrices, $a = \text{rank}(\mathbf{C}) \leq q$ and $b = \text{rank}(\mathbf{U}) \leq p$. Subtraction of a constant matrix Θ_0 has no effect on estimation and precision and we may ignore Θ_0 without loss of generality when discussing nearly all estimation properties.

An alternative form is $\tau = \mathbf{L}\beta_s = \mathbf{L}\text{vec}(\mathbf{B})$, with \mathbf{L} constant. For convenience we assume $\ell = \text{rank}(\mathbf{L}) \leq pq$. An additive constant may be ignored without loss of generality. Although the $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ form is more common, the $\tau = \mathbf{L}\beta_s$ form is also valid and more general in that some secondary parameters can be defined in the $\mathbf{L}\beta_s$ form but cannot be defined in the $\mathbf{C}\mathbf{B}\mathbf{U}$ form.

One can easily verify the above statements as follows. Each element of $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ is a linear combination of the elements of \mathbf{B} and, therefore, of the elements of β_s . Any linear combinations of the elements of \mathbf{B} can be written in the form $\mathbf{L}\beta_s$. It is easy to find counterexamples to the converse.

Theorem 12.7 Having $\mathbf{B} \in \mathbb{R}^{q \times p}$ allows defining $\beta_s = \text{vec}(\mathbf{B})$ and considering fixed constants \mathbf{C} , \mathbf{U} , and \mathbf{L} .

- (a) If $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ and $\tau = \mathbf{L}\beta_s$, then $\text{vec}(\Theta) = \tau \forall \mathbf{B} \Leftrightarrow \mathbf{L} = (\mathbf{U}' \otimes \mathbf{C})$.
- (b) If $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ and $\tau = \text{vec}(\Theta)$, then \mathbf{L} exists such that $\tau = \mathbf{L}\beta_s \forall \mathbf{B}$.
- (c) If $\tau = \mathbf{L}\beta_s$, then $\{\mathbf{C}, \mathbf{U}\}$ may or may not exist such that $\tau = \text{vec}(\mathbf{C}\mathbf{B}\mathbf{U}) \forall \mathbf{B}$.

Proof. (a) follows directly from the properties of $\text{vec}()$ and $\mathbf{A} \otimes \mathbf{B}$.

For (b), choosing $\mathbf{L} = (\mathbf{U}' \otimes \mathbf{C})$ satisfies the claim.

(c) A counterexample is $\mathbf{L} = [1 \ 0 \ 0 \ -1]$, $\mathbf{B} = \begin{bmatrix} \beta_{11} & \beta_{12} \\ \beta_{21} & \beta_{22} \end{bmatrix}$, $\beta_s = [\beta_{11} \ \beta_{21} \ \beta_{12} \ \beta_{22}]'$, and $\tau = \mathbf{L}\beta_s = (\beta_{11} - \beta_{22})$. If $1 \times q$ \mathbf{C} and $p \times 1$ \mathbf{U} exist such that $\tau = \mathbf{L}\beta_s = \mathbf{C}\mathbf{B}\mathbf{U}$, then $\mathbf{C}\mathbf{B}\mathbf{U} = \sum_{i=1}^2 \sum_{j=1}^2 C_i \beta_{ij} U_j = (\beta_{11} - \beta_{22})$ must hold $\forall \mathbf{B}$. The last equation requires the impossible, that (a) some elements of $\{\mathbf{C}, \mathbf{U}\}$ are zero, and (b) none of the elements of $\{\mathbf{C}, \mathbf{U}\}$ are zero. The supposition that $\{\mathbf{C}, \mathbf{U}\}$ exists has led to a contradiction. □

Lemma 12.6 If θ ($a \times 1$) is an estimable secondary (or primary) parameter for a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ and $\hat{\theta}$ is the BLUE of θ , then a tertiary (or secondary) parameter vector is $\gamma = \mathbf{T}\theta$, with $g = \text{rank}(\mathbf{T}) \leq a$. The BLUE of γ is $\hat{\gamma} = \mathbf{T}\hat{\theta}$.

Proof. If $\hat{\theta} = \text{BLUE}(\theta)$, then (1) $\hat{\theta}$ is an LUE, which ensures $\exists \mathbf{A}$ such that $\hat{\theta} = \mathbf{A}\mathbf{y}$ and $\text{E}(\mathbf{A}\mathbf{y}) = \theta$. Hence $\hat{\gamma}$ is a LUE since $\text{E}(\hat{\gamma}) = \text{E}[(\mathbf{T}\mathbf{A})\mathbf{y}] = \mathbf{T}\text{E}(\mathbf{A}\mathbf{y}) = \mathbf{T}\theta = \gamma$. Also, $\hat{\theta} = \text{BLUE}(\theta)$ implies (2) $\hat{\theta}$ is the “best” LUE. For any other LUE, say, $\hat{\theta}_1 = \mathbf{A}_1\mathbf{y}$ ($\neq \mathbf{A}$), $\mathcal{V}(\mathbf{s}'\mathbf{A}\mathbf{y}) = \mathbf{s}'\mathcal{V}(\mathbf{A}\mathbf{y})\mathbf{s}$ and $\mathcal{V}(\mathbf{s}'\mathbf{A}_1\mathbf{y}) = \mathbf{s}'\mathcal{V}(\mathbf{A}_1\mathbf{y})\mathbf{s}$. By Theorem 11.9, $\mathcal{V}(\mathbf{A}_1\mathbf{y}) - \mathcal{V}(\mathbf{A}\mathbf{y})$ is positive semidefinite. The definition of “positive semidefinite” immediately implies $\mathcal{V}(\mathbf{s}'\mathbf{A}_1\mathbf{y}) \geq \mathcal{V}(\mathbf{s}'\mathbf{A}\mathbf{y}) \forall \mathbf{s} \in \mathbb{R}^a$. For $\mathbf{s} \in \{\mathbf{s} : \mathbf{s}' = \mathbf{k}'\mathbf{T} \text{ with } \mathbf{k} \in \mathbb{R}^g\} \subseteq \mathbb{R}^a$ it follows that $\mathcal{V}(\mathbf{s}'\mathbf{A}_1\mathbf{y}) \geq \mathcal{V}(\mathbf{s}'\mathbf{A}\mathbf{y})$ for all such \mathbf{s} and $\mathcal{V}(\mathbf{k}'\mathbf{T}\mathbf{A}_1\mathbf{y}) \geq \mathcal{V}(\mathbf{k}'\mathbf{T}\mathbf{A}\mathbf{y}) \forall \mathbf{k} \in \mathbb{R}^g$. \square

Theorem 12.8 For $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$, $\beta_s = \text{vec}(\mathbf{B})$, $r = \text{rank}(\mathbf{X}) \leq q$, and estimable secondary parameter $\tau = \mathbf{L}\beta_s$ with $\tau \ell \times 1$, $\ell = \text{rank}(\mathbf{L}) \leq pq$, if $\hat{\beta}_s$ is the BLUE of β_s , then $r = q$ and the BLUE of τ is $\hat{\tau} = \mathbf{L}\hat{\beta}_s$.

Proof. The theorem follows from the form of the likelihood and from a univariate estimation theorem. The details are left to the reader.

Corollary 12.8 For estimable secondary parameter matrix $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ with $a = \text{rank}(\mathbf{C}) \leq q$ and $b = \text{rank}(\mathbf{U}) \leq p$, the BLUE is $\hat{\Theta} = \mathbf{C}\hat{\mathbf{B}}\mathbf{U}$.

Proof. The result follows immediately from Theorem 11.14 and the preceding theorems in the present chapter. The details are left to the reader as an exercise.

Theorem 12.9 (a) For $\text{GLM}_{N,p,q}\text{LTFR}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and two-condition inverse $(\mathbf{X}'\mathbf{X})^-$,

$$\tilde{\mathbf{B}} \sim \mathcal{SN}_{p,q}[(\mathbf{X}'\mathbf{X})^-(\mathbf{X}'\mathbf{X})\mathbf{B}, (\mathbf{X}'\mathbf{X})^-, \Sigma]. \quad (12.41)$$

(b) A $\text{GLM}_{N,p,q}\text{FR}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ has

$$\hat{\mathbf{B}} \sim \mathcal{N}_{p,q}[\mathbf{B}, (\mathbf{X}'\mathbf{X})^{-1}, \Sigma]. \quad (12.42)$$

(c) For $\text{rank}(\mathbf{X}) \leq q$ and any one-condition inverse, estimable $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ gives

$$\hat{\Theta} - \Theta_0 \sim (\mathcal{S})\mathcal{N}_a[\Theta - \Theta_0, \mathbf{C}(\mathbf{X}'\mathbf{X})^-\mathbf{C}', \mathbf{U}'\Sigma\mathbf{U}]. \quad (12.43)$$

Proof. Left as an exercise.

12.5 ESTIMATION WITH MULTIVARIATE RESTRICTIONS

Writing $GLM_{N,p,q}(Y_i; X_i B | R_x B R_y = A, \Sigma)$ indicates an *explicitly* restricted multivariate model. The $(q \times p)$ matrix B has between-ISU (independent sampling unit, “subject”) restrictions defined by known and constant matrices R_x of dimension $a \times q$ with $a \leq q$, $\text{rank}(R_x) = a$. The $(q \times p)$ matrix B has within-subject restrictions defined by R_y of dimension $p \times b$ with $p \geq b$ and $\text{rank}(R_y) = b$. The $a \times b$ constant A plays a role in both sorts of restrictions. The matrices R_x , R_y , and A must all be chosen without knowledge of the data, prior to data collection. As in univariate models, matrix R_x implicitly specifies restrictions about X , the between-subjects design matrix. The matrix R_y implicitly specifies restrictions about Y and can be thought of as indirectly defining a within-subject design matrix. Naturally all results about restrictions in univariate models occur in the special case of the multivariate model with $p = 1$.

Example 12.2 A set of *ipsative* measures add to a known constant. Such variables occur naturally in the study of allocation of behavior and many other areas. The pandemic of obesity in the United States has led scientists to assess how much time Americans allocate to sleeping, eating, watching television, walking, etc. Each person's allocations always add to 24 hours. The implicit constraint creates a singular covariance matrix (Σ), and difficulties in defining scientifically appealing parameter estimators which account for the constraint.

Example 12.3 Pan (2003) evaluated models of land use allocation among rural farmers in the Amazon Basin of South America. With farm as the independent sampling unit, the multivariate response of interest was % land in crops, % land in pasture, % land fallow. The formulation corresponds to the restrictions $B R_y = A$, with $R_y = [1 \ 1 \ 1]$ and $A = [100]$.

Example 12.4 The same constraint arises in the study of diets as predictors of cancer. Nutritional epidemiologists seek to build models of % calories from fat, % calories from protein, % calories from carbohydrates, and % calories from alcohol.

Having restrictions only between subjects requires $R_y = I_p$ and simplifies the model statement of an explicitly restricted linear model to the form $GLM_{N,p,q}(Y_i; X_i B | R_x B = A, \Sigma)$. Any such model can be transformed to a linearly equivalent *unrestricted* linear model $GLM_{N,p,q}(Y_{iu}; X_{iu} B_u, \Sigma_u)$ with all dimensions the same. Adding within-subject restrictions corresponds to choosing $R_y \neq I_p$ and writing $GLM_{N,p,q}(Y_i; X_i B | R_x B R_y = A, \Sigma)$.

The same statement does not hold, except in special cases, for a model having within-subject restrictions, $GLM_{N,p,q}(Y_i; X_i B | R_x B R_y = A, \Sigma)$. Without loss of generality, for the purposes of the present argument, we may assume $R_x = I_q$ and discuss $GLM_{N,p,q}(Y_i; X_i B | B R_y = A, \Sigma)$. A linearly equivalent unrestricted model can be guaranteed to exist in such a model only if $b = p = \text{rank}(R_y)$.

However, in practice, the restriction statement often implicitly defines the only parameters of interest. If so, analysis of the (within-subject) restricted model only loses access to scientifically irrelevant parameters. Equally importantly, and not surprisingly, properties of the restricted model suffice to fully specify the distributions of estimators and tests available in the restricted model.

The next theorem has many uses. The results provide ways of defining linearly equivalent multivariate GLMs and GGLMs in the presence of restrictions. Looking ahead, testing a general linear hypothesis may be cast as specifying a set of restrictions. The various forms help simplify derivations of estimators, tests, and the associated distributions.

Example 12.5 A paired data t test can be cast as a test in a multivariate linear model with $p = 2$ or as a univariate model with $p = 1$ and the responses being difference scores. In both cases, $\mathbf{X} = \mathbf{1}_N$. The following theorem allows expressing the $p = 2$ model as a restriction of the $p = 1$ model. The restriction matrix is $\mathbf{U} = [1 \ -1]'$.

Theorem 12.10 Model 1 is $\text{GLM}_{N,p,q}(Y_i; X_i B | R_x B R_y = A, \Sigma)$, with R_x ($a \times q$) rank $a \leq q$ and R_y ($p \times b$) rank $b \leq p$. If $R_{y\perp}$ is any $p \times (p - b)$ matrix such that $\mathbf{T} = [R_y \ R_{y\perp}]$ is $p \times p$ and full rank, it defines

$$\Gamma = \mathbf{B}\mathbf{T} = [BR_y \ BR_{y\perp}] = \begin{bmatrix} \Gamma_1 & \Gamma_2 \\ q \times b & q \times (p-b) \end{bmatrix}. \tag{12.44}$$

Here $\mathbf{B} = [\Gamma_1 \ \Gamma_2]\mathbf{T}^{-1}$. If $R_x \Gamma_1 = \mathbf{A}$ is consistent, then Γ_1 satisfies the restrictions if and only if, for some $q \times b \Delta$ (with sign adjustable),

$$\Gamma_1 = R_x^- \mathbf{A} \pm (R_x^- R_x - I_q) \Delta. \tag{12.45}$$

In turn, $R_x \Gamma_1 = \mathbf{A}$ is equivalent to the original restrictions, and models 2, 3, 4, and 5, defined below, satisfy the restrictions. Models 1, 2, 4, and 5 are always linearly equivalent to each other, and linearly equivalent to model 3 if R_y is square and full rank.

Model 2 is $\text{GLM}_{N,p,q}(Y_i \mathbf{T}, X_i [\Gamma_1 \ \Gamma_2] | R_x \Gamma_1 = \mathbf{A}, \mathbf{T}' \Sigma \mathbf{T})$.

Model 3 is $\text{GLM}_{N,b,q}(Y_{i3}, X_{i3} \Delta, R_y' \Sigma R_y)$, with $Y_3 = Y R_y - X R_x^- \mathbf{A}$ and $X_3 = X (R_x^- R_x - I_q)$ and ignores Γ_2 . If $b = p = \text{rank}(R_y)$, it is linearly equivalent.

Model 4 is $\text{GGLM}_{N,p,q}(\mathbf{y}_s; X_s \beta_s | R_s \beta_s = \mathbf{a}_s, \Sigma \otimes I_N)$, with $\mathbf{y}_s = \text{vec}(Y)$, $X_s = (I_p \otimes X)$, $\beta_s = \text{vec}(B)$, R_s ($ab \times q$), \mathbf{a}_s ($ab \times 1$), $R_s = (\mathbf{T}' \otimes R)$, and $\mathbf{a}_s = \text{vec}(A)$.

Model 5 is $\text{GGLM}_{N,p,q}(\mathbf{u}; Z \gamma_s, \Sigma \otimes I_N)$, with $\mathbf{u} = \mathbf{y}_s - X_s R_s^- \mathbf{a}_s$, $Z = X_s (R_s^- R_s - I_{qp})$, γ_s of dimension $qp \times 1$, and Z of dimension $Np \times qp$.

Proof. Expressions for Γ_1 in terms of Δ may be derived in terms of an arbitrary column j of Γ_1 , say γ_{1j} ($q \times 1$). Here $R_x \gamma_{1j} = \mathbf{a}_i$ ($a \times 1$) $\Leftrightarrow \gamma_{1j} = R_x^- \mathbf{A}_j + (R_x^- R_x - I_q) \delta_j$ for some δ_j ($q \times 1$). Results on univariate models ensure the desired properties hold. The remainder of the proof is also based on univariate results and is left to the reader as an exercise. \square

**12.6 UNRESTRICTED ESTIMATION WITH COMPOUND SYMMETRY:
THE “UNIVARIATE” APPROACH TO REPEATED MEASURES**

As detailed in Lemma 1.33, a $p \times p$ compound symmetric covariance matrix may be written

$$\begin{aligned} \Sigma &= \sigma^2[\mathbf{1}_p \mathbf{1}'_p \rho + \mathbf{I}_p(1 - \rho)] \\ &= \mathbf{V} \text{Dg}(\lambda) \mathbf{V}' \\ &= [\mathbf{v}_B \quad \mathbf{V}_W] \begin{bmatrix} \lambda_1 & \mathbf{0} \\ \mathbf{0} & \lambda_2 \mathbf{I}_{p-1} \end{bmatrix} \begin{bmatrix} \mathbf{v}'_B \\ \mathbf{V}'_W \end{bmatrix}. \end{aligned} \tag{12.46}$$

The $p \times p$ matrix $\text{Dg}(\lambda) = \text{Dg}(\lambda_1, \lambda_2 \mathbf{1}_{p-1})$ has $\lambda_1 = \sigma^2[1 + (p - 1)\rho]$ and $\lambda_2 = \sigma^2(1 - \rho)$, with known eigenvectors.

Theorem 12.11 Model 1, $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$, has positive definite compound symmetric covariance $\Sigma = \mathbf{V} \text{Dg}(\lambda) \mathbf{V}'$. The eigenvectors provide a *known, constant and exact* transformation to a model with uncorrelated observations, model 2, $\text{GLM}_{N,p,q}[\mathbf{Y}_i \mathbf{V}; \mathbf{X}_i \mathbf{B} \mathbf{V}, \text{Dg}(\lambda)]$. With Gaussian data the columns of $\mathbf{Y} \mathbf{V}$, as well as the rows, are independent. Column 1 has variance $\lambda_1 = \sigma^2[1 + (p - 1)\rho]$ and the other $p - 1$ columns have variance $\lambda_2 = \sigma^2(1 - \rho)$.

Proof. Lemma 1.33 provides most results needed. Model transformation gives

$$\begin{aligned} \mathbf{Y} &= \mathbf{X} \mathbf{B} + \mathbf{E} \\ \mathbf{Y} \mathbf{V} &= \mathbf{X} \mathbf{B} \mathbf{V} + \mathbf{E} \mathbf{V}, \end{aligned} \tag{12.47}$$

$E[\text{row}_i(\mathbf{E} \mathbf{V})] = \{E[\text{row}_i(\mathbf{E})]\} \mathbf{V} = \mathbf{0}$, and $\mathcal{V}\{\text{row}_i(\mathbf{E} \mathbf{V})\}' = \mathbf{V}' \Sigma \mathbf{V} = \text{Dg}(\lambda)$. \square

Corollary 12.11 If the original data are Gaussian, then all observations split into one set of N observations which exactly follow model 2B (between), a *univariate* $\text{GLM}_{N,q}()$, and a second set of $N(p - 1)$ observations which exactly follow model 2W (within), a *univariate* $\text{GLM}_{N(p-1),q}()$.

Proof. If $\mathbf{y}_B = \mathbf{Y} \mathbf{v}_B$ ($N \times 1$) and $\mathbf{Y}_W = \mathbf{Y} \mathbf{V}_W$ [$N \times (p - 1)$], then

$$\begin{aligned} \mathbf{Y} \begin{bmatrix} \mathbf{v}_B & \mathbf{V}_W \end{bmatrix} &= \mathbf{X} \mathbf{B} \begin{bmatrix} \mathbf{v}_B & \mathbf{V}_W \end{bmatrix} + \mathbf{E} \begin{bmatrix} \mathbf{v}_B & \mathbf{V}_W \end{bmatrix} \\ \begin{bmatrix} \mathbf{y}_B & \mathbf{Y}_W \end{bmatrix} &= \mathbf{X} \begin{bmatrix} \mathbf{B} \mathbf{v}_B & \mathbf{B} \mathbf{V}_W \end{bmatrix} + \begin{bmatrix} \mathbf{e}_B & \mathbf{E}_W \end{bmatrix}. \end{aligned} \tag{12.48}$$

If $\boldsymbol{\beta}_B = \mathbf{B} \mathbf{v}_B$ ($q \times 1$), extracting the model (2B) for the first column gives

$$\mathbf{y}_B = \mathbf{X} \boldsymbol{\beta}_B + \mathbf{e}_B. \tag{12.49}$$

Furthermore

$$\mathbf{e}_B \sim \mathcal{N}_N(\mathbf{0}, \lambda_1 \mathbf{I}_N) \tag{12.50}$$

$$\mathbf{y}_B \sim \mathcal{N}_N(\mathbf{X} \boldsymbol{\beta}_B, \lambda_1 \mathbf{I}_N). \tag{12.51}$$

Hence model 2B meets the assumptions of $\text{GLM}_{N,q}(y_{Bi}; \mathbf{X}_i \boldsymbol{\beta}_B, \lambda_1)$ with Gaussian errors. With similar notation, $\mathbf{B}_W = \mathbf{B} \mathbf{V}_W$ is $q \times (p - 1)$ and

$$\mathbf{Y}_W = \mathbf{X}\mathbf{B}_W + \mathbf{E}_W. \quad (12.52)$$

Here

$$\begin{aligned} \mathbf{E}_W &\sim \mathcal{N}_{N,(p-1)}(\mathbf{0}, \mathbf{I}_N, \lambda_2 \mathbf{I}_{p-1}) \Leftrightarrow \\ \mathbf{e}_W = \text{vec}(\mathbf{E}_W) &\sim \mathcal{N}_{N(p-1)}(\mathbf{0}, \lambda_2 \mathbf{I}_{N(p-1)}), \end{aligned} \quad (12.53)$$

$$\begin{aligned} \mathbf{Y}_W &\sim \mathcal{N}_{N,(p-1)}(\mathbf{X}\mathbf{B}_W, \mathbf{I}_N, \lambda_2 \mathbf{I}_{p-1}) \Leftrightarrow \\ \mathbf{y}_W = \text{vec}(\mathbf{Y}_W) &\sim \mathcal{N}_{N(p-1)}[\text{vec}(\mathbf{X}\mathbf{B}_W), \lambda_2 \mathbf{I}_{N(p-1)}]. \end{aligned} \quad (12.54)$$

Observing

$$\begin{aligned} \text{vec}(\mathbf{X}\mathbf{B}_W) &= (\mathbf{I}_{p-1} \otimes \mathbf{X})\text{vec}(\mathbf{B}_W) \\ &= \mathbf{X}_W \boldsymbol{\beta}_W \end{aligned} \quad (12.55)$$

allows writing

$$\begin{aligned} \text{vec}(\mathbf{Y}_W) &= (\mathbf{I}_{p-1} \otimes \mathbf{X})\text{vec}(\mathbf{B}_W) + \text{vec}(\mathbf{E}_W) \Leftrightarrow \\ \mathbf{y}_W &= \mathbf{X}_W \boldsymbol{\beta}_W + \mathbf{e}_W. \end{aligned} \quad (12.56)$$

Finally, the last equation satisfies univariate $\text{GLM}_{N(p-1),q}[y_{W,i}; \text{row}_i(\mathbf{X}_W)\boldsymbol{\beta}_W, \lambda_2]$ with Gaussian errors, model 2W. \square

Model 2B describes differences between subjects (ISUs), while model 2W describes differences within subjects. The formulation provides a basis of the theory underlying the “univariate” approach to repeated measures and the “uncorrected” UNIREP test. Muller and Barton (1989) included many details.

EXERCISES

Here \mathbf{e}_1 is a vector with a 1 as the first element and zeros as all other elements: $\mathbf{e}_1 = [1 \ 0 \ \cdots \ 0]'$. In different exercises \mathbf{e}_1 can have different dimensions.

True/False Exercises. Several propositions are stated below. For each choose one of the following: T if the proposition is true and F if the proposition is false. Be careful! Some propositions may be tricky. Most are either correct theorems or a slight modification of correct theorems (in which case the proposition is false).

For each proposition you mark as “false,” give either (1) a brief counterexample or (2) a brief remark on why the proposition fails to be true. This should make clear which aspects of the proposition you reject. In addition to intentional errors in the propositions, there is always the possibility of unintentional typographical errors. Therefore, your brief remarks are important. If you mark a proposition “true” when it is in fact false, it will be inferred that your proof would contain one or more errors. Therefore, do not attach proofs or comments for propositions marked true.

The following assumption and definitions apply to exercises 12.1–12.6. Assume $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and $\text{rank}(\mathbf{X}) = r \leq q$. Also $\widehat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} = [\widehat{\beta}_1 \ \widehat{\beta}_2 \ \cdots \ \widehat{\beta}_p]$, $\widehat{\Sigma} = \mathbf{Y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y}/(N - r)$, $\widehat{\beta}_1 = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{e}_1$, and $\widehat{\sigma}_1^2 = \mathbf{e}_1'\mathbf{Y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y}\mathbf{e}_1/(N - r)$.

T F 12.1 Proposition: $\widehat{\beta}_1 \sim \mathcal{N}_q[\mathbf{C}\mathbf{B}\mathbf{e}_1, \mathbf{e}_1'\Sigma\mathbf{e}_1(\mathbf{X}'\mathbf{X})^{-1}]$ of rank $r \leq q$.

T F 12.2 Proposition: $\widehat{\sigma}_1^2/(\mathbf{e}_1'\Sigma\mathbf{e}_1) \sim \chi^2(N - r, 0)$.

T F 12.3 Proposition: $\widehat{\beta}_1$ and $\widehat{\sigma}_1^2$ are uncorrelated.

T F 12.4 Additionally assume $\boldsymbol{\theta} = \mathbf{C}\mathbf{B}\mathbf{e}_1$ and $\widehat{\boldsymbol{\theta}} = \mathbf{C}\widehat{\mathbf{B}}\mathbf{e}_1$ are $a \times 1$, $a \leq q$.

Proposition: $\widehat{\boldsymbol{\theta}} \sim \mathcal{N}_a[\mathbf{C}\mathbf{B}\mathbf{e}_1, \mathbf{e}_1'\Sigma\mathbf{e}_1\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']$ and $\text{rank}[\mathcal{V}(\widehat{\boldsymbol{\theta}})] = a$.

The following assumptions and notation apply to exercises 12.5–12.8. Assume $\text{GLM}_{N,p,q}\text{FR}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors, $N \gg q$, and $\text{rank}(\Sigma) = p$. Also $\widehat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}$, $\mathbf{S} = \mathbf{Y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y}/(N - q)$, $\mathbf{y}_s = \text{vec}(\mathbf{Y})$, $\boldsymbol{\beta}_s = \text{vec}(\mathbf{B})$, and $\widehat{\boldsymbol{\beta}}_s = \text{vec}(\widehat{\mathbf{B}})$.

T F 12.5 Proposition: $\widehat{\boldsymbol{\beta}}_s \sim \mathcal{N}_{pq}[\boldsymbol{\beta}_s, \Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1}]$ and $\text{rank}[\mathcal{V}(\widehat{\boldsymbol{\beta}}_s)] = pq$.

T F 12.6 Proposition: $\mathbf{S} \sim \mathcal{W}_p(N - q, \Sigma, \mathbf{0})$.

T F 12.7 Proposition: $\widehat{\mathbf{B}}$ and \mathbf{S} are independent.

T F 12.8 Proposition: $E(\mathbf{S}) = \Sigma$ and $\mathcal{V}(s_{ij}) = (N - q)\sigma_{ij}^2 + (N - q)^2\sigma_{ii}\sigma_{jj}$.

12.9 Three models may be written for the same data as follows: For model 1

$$\begin{aligned} \mathbf{Y} &= (\mathbf{1}_n \otimes \mathbf{I}_3)\mathbf{B}_1 + \mathbf{E} \\ \mathbf{Y} &= \mathbf{X}_1\mathbf{B}_1 + \mathbf{E}. \end{aligned}$$

For model 2

$$\begin{aligned} \mathbf{Y} &= \left(\mathbf{1}_n \otimes \begin{bmatrix} 1 & 0 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \right) \mathbf{B}_2 + \mathbf{E} \\ \mathbf{Y} &= \mathbf{X}_2\mathbf{B}_2 + \mathbf{E}. \end{aligned}$$

For model 3

$$\begin{aligned} \mathbf{Y} &= \left(\mathbf{1}_n \otimes \begin{bmatrix} 1 & 1 & 0 & 0 & -1 & -1 \\ 1 & 0 & 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 & 0 & 1 \end{bmatrix} \right) \mathbf{B}_3 + \mathbf{E} \\ \mathbf{Y} &= \mathbf{X}_3\mathbf{B}_3 + \mathbf{E}. \end{aligned}$$

Each meets the assumptions of a $\text{GLM}_{N,p,qk}(\mathbf{Y}_i; \mathbf{X}_{i,k}\mathbf{B}_k, \Sigma)$.

12.9.1 Prove that model 1 and model 2 are linearly equivalent.

12.9.2 Prove that model 1 and model 3 are linearly equivalent.

12.9.3 matrices $\mathbf{R}_x = [\mathbf{0} \ \mathbf{I}_3 \ \mathbf{0}]$ and $\mathbf{A} = \mathbf{0}$ specify model 4 $\text{GLM}_{N,p,q_3}(\mathbf{Y}_i; \mathbf{X}_{i,3}\mathbf{B}_3 | \mathbf{R}_x\mathbf{B}_3 = \mathbf{A}, \Sigma)$ as a restriction of model 3. Specify the dimensions of \mathbf{A} and submatrices in \mathbf{R}_x .

12.9.4 Find a full-rank and unrestricted model 5, equivalent to model 4, with the

additional restriction that all elements of \mathbf{B}_5 are also elements of \mathbf{B}_3 . This may be done by inspection (no proof needed).

CHAPTER 13

Estimation for Generalizations of Multivariate Models

13.1 MOTIVATION

As detailed in the preceding two chapters, the “least squares” assumptions of the univariate and multivariate GLM allow finding unbiased estimators of the primary parameters and broad classes of secondary parameters. The additional assumption of Gaussian errors leads to the estimators also satisfying likelihood principles. In all cases the estimates can be computed in closed form, with a noniterative algorithm for solving a system of simultaneous linear equations.

Do more general linear models allow such nice results? For a range of generalizations of the multivariate linear model, the answer is essentially “Yes.” Although usually much less convenient to compute, estimates satisfying generalized least squares or likelihood criteria can be found. Furthermore, the expected value parameter estimators are usually unbiased or nearly so. Covariance parameter estimators can range from unbiased to substantially biased.

As discussed in Chapter 4, models for growth curves, seemingly unrelated regressions, multiple design matrix settings, and doubly multivariate settings can be cast as generalizations of the GLM. The references given there contain details about particular techniques.

13.2 CRITERIA AND ALGORITHMS

Given Gaussian errors, ordinary least squares (OLS) estimators coincide with the likelihood estimators in the univariate and multivariate GLM. The same holds true for *exact* weighted least squares (WLS), as seen in results in the previous two chapters for the GGLM.

Definition 13.1 (a) Using estimated weights (covariance matrix) in weighted least squares (WLS) formulas defines *approximate* weighted least squares (AWLS) estimators.

(b) Using WLS formulas to iteratively estimate expected values and covariance defines *iterative approximate least squares* (ITAWLS).

AWLS requires an initial computation of an OLS estimate for the expected-value parameters, followed by computation of covariance parameter estimates in terms of the residuals. Except in very special cases (e.g., Gaussian errors in special patterns), the two-step method does not produce maximum likelihood estimates.

Although much more efficient algorithms can often be found, iterating the two-step AWLS method (ITAWLS) yields maximum likelihood estimates for a variety of models with Gaussian errors. The covariance estimate from the second step allows computing an updated estimate of the expected-value parameters. In turn, a new covariance estimate may be computed, etc.

13.3 WEIGHTED ESTIMATION OF B AND Σ

Weighted estimation in a multivariate model implies considering $Y \sim \mathcal{N}_{N,p}(XB, D, \Sigma)$, with *known* X , full-rank $D = F_D F_D'$, and unknown B and Σ . Given the assumptions, $\text{vec}(Y') \sim \mathcal{N}_{Np}[(X \otimes I_p)\text{vec}(B'), D \otimes \Sigma]$. Multiplication by F_D^{-1} gives $Y_D = F_D^{-1}Y \sim \mathcal{N}_{N,p}(F_D^{-1}XB, I, \Sigma)$. If $X_D = F_D^{-1}X$, then Y_D corresponds to $\text{GLM}_{N,p,q}[\text{row}_i(Y_D), \text{row}_i(X_D)B, \Sigma] \Leftrightarrow \text{GLM}_{N,p,q}(F_D^{-1}Y_i, F_D^{-1}X_iB, \Sigma)$.

Theorem 13.1 Matrices Y ($N \times p$), X ($N \times q$), $\text{rank}(X) \leq q$, positive definite $D = D'$ ($N \times N$), and positive definite or positive semidefinite $\Sigma = \Sigma'$ ($p \times p$) define $\text{GGLM}_{Np,q}[\text{vec}(Y'); (X \otimes I_p)\text{vec}(B'), D \otimes \Sigma]$, with Gaussian errors. The same assumptions give $\text{GGLM}_{Np,q}[\text{vec}(Y'L); (L'X \otimes I_p)\text{vec}(B'), I_N \otimes \Sigma]$. The latter model has a corresponding multivariate GLM satisfying the least squares assumptions with Gaussian errors and $\text{rank}(L'X) = \text{rank}(X)$.

Proof. With $Y_i = \text{row}_i(Y)$ and Y' $p \times N$, $\text{vec}(Y')$ is the “vertical concatenation of rows.” Hence

$$\begin{aligned} \mathcal{V}[\text{vec}(Y')] &= D \otimes \Sigma \\ &= \begin{bmatrix} d_{11}\Sigma & d_{12}\Sigma & \cdots & d_{1N}\Sigma \\ d_{21}\Sigma & d_{22}\Sigma & \cdots & d_{2N}\Sigma \\ \vdots & \vdots & \ddots & \vdots \\ d_{N1}\Sigma & d_{N2}\Sigma & \cdots & d_{NN}\Sigma \end{bmatrix}, \end{aligned} \tag{13.1}$$

and $\text{vec}(L'Y) = (I_p \otimes L')\text{vec}(Y)$ has covariance $(I_p \otimes L')(\Sigma \otimes D)(I_p \otimes L) = \Sigma \otimes I_N$, while $\text{vec}(Y'L) = (L' \otimes I_p)\text{vec}(Y')$ has covariance $(L' \otimes I_p)(D \otimes \Sigma)(L \otimes I_p) = I_N \otimes \Sigma$ and mean $(L' \otimes I_p)(X \otimes I_p)\text{vec}(B') = (L'X \otimes I_p)\text{vec}(B')$. The proof is completed by appealing to properties of a direct-product matrix Gaussian. \square

Theorem 13.2 For $\text{GGLM}_{N,p,q}[\text{vec}(\mathbf{Y}'); (\mathbf{X} \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}'), \mathbf{D} \otimes \mathbf{\Sigma}]$ with Gaussian errors, $r = \text{rank}(\mathbf{X}) \leq q$, $\mathbf{\Theta} = \mathbf{C}\mathbf{B}$ is estimable, unknown $\mathbf{\Sigma} = \mathbf{\Sigma}'$ ($p \times p$) is positive definite, while $\mathbf{D} = \mathbf{D}'$ ($N \times N$) is known and positive definite. The joint MLEs of \mathbf{B} , $\mathbf{\Sigma}$, and $\mathbf{\Theta}$ are

$$\tilde{\mathbf{B}} = (\mathbf{X}'\mathbf{D}^{-1}\mathbf{X})^{-1}\mathbf{X}'\mathbf{D}^{-1}\mathbf{Y} \tag{13.2}$$

$$\tilde{\mathbf{\Sigma}} = (\mathbf{Y} - \mathbf{X}\tilde{\mathbf{B}})'\mathbf{D}^{-1}(\mathbf{Y} - \mathbf{X}\tilde{\mathbf{B}})/N \tag{13.3}$$

$$\hat{\mathbf{\Theta}} = \mathbf{C}\tilde{\mathbf{B}}. \tag{13.4}$$

Both $\tilde{\mathbf{\Sigma}}$ and $\hat{\mathbf{\Theta}}$ are invariant to $\tilde{\mathbf{B}}$, which is unique if $r = q$. The unbiased restricted maximum likelihood estimator $\hat{\mathbf{\Sigma}} = \tilde{\mathbf{\Sigma}}N/(N - r)$ is typically used.

Proof. The proof is left as an exercise.

13.4 TRANSFORMATIONS AMONG GROWTH CURVE DESIGNS

Often a $\text{GCM}_{N,p,q,m}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}\mathbf{T}, \mathbf{\Sigma})$ leads to interest in a reduced model for the within-subject dimension. As in all our discussions of the GCM, for convenience of exposition we assume that columns of \mathbf{Y}_i differ (only) in regard to the time a repeated measure was collected. With polynomial coding, the reduced (within-subject) model can typically be created from the full (within-subject) model by deleting higher-order terms corresponding to trailing columns of \mathbf{B} . If so, the expected value for the full model may be written

$$\begin{aligned} E(\mathbf{Y}|\mathbf{X}, \mathbf{T}) &= \mathbf{X}\mathbf{B}\mathbf{T} \\ &= \mathbf{X}[\mathbf{B}_1 \ \mathbf{B}_2] \begin{bmatrix} \mathbf{T}_1 \\ \mathbf{T}_2 \end{bmatrix} \\ &= \mathbf{X}\mathbf{B}_1\mathbf{T}_1 + \mathbf{X}\mathbf{B}_2\mathbf{T}_2, \end{aligned} \tag{13.5}$$

with corresponding expected value for the reduced model of

$$\begin{aligned} E(\mathbf{Y}|\mathbf{X}, \mathbf{T}, \mathbf{B}_2 = \mathbf{0}) &= \mathbf{X}[\mathbf{B}_1 \ \mathbf{0}] \begin{bmatrix} \mathbf{T}_1 \\ \mathbf{T}_2 \end{bmatrix} \\ &= \mathbf{X}\mathbf{B}_1\mathbf{T}_1. \end{aligned} \tag{13.6}$$

With $m \times p$ \mathbf{T} for $1 \leq m \leq p$ and $m_1 \times p$ \mathbf{T}_1 , necessarily \mathbf{T}_2 must be $(m - m_1) \times p$ for $1 \leq m_1 \leq m \leq p$. The fully saturated (within-subject) model has $m = p$ and $\mathbf{T} = \{d_j^{i-1}\}$. Equivalently, for ISU i the expected values are written $E(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{T}) = \mathbf{X}_i\mathbf{B}\mathbf{T} = \mathbf{X}_i\mathbf{B}_1\mathbf{T}_1 + \mathbf{X}_i\mathbf{B}_2\mathbf{T}_2$ and $E(\mathbf{Y}_i|\mathbf{X}_i, \mathbf{T}, \mathbf{B}_2 = \mathbf{0}) = \mathbf{X}_i\mathbf{B}_1\mathbf{T}_1$.

A similar construction relating a full and reduced model for between-subject design applies when the columns of \mathbf{X} contain polynomials for a continuous predictor. A scientist using natural polynomial coding and a participant's age as the basic predictor can use $\mathbf{X}_i = \{x_{ij}\}$ with $x_{ij} = \text{age}_i^{j-1}$, giving $\mathbf{X}_i = [1 \ \text{age}_i \ \text{age}_i^2 \ \cdots \ \text{age}_i^{q-1}]$. The strong constraints relating the univariate GLM to all variations of the multivariate GLM, including the GCM, guarantee

univariate techniques for coding \mathbf{X} and reducing models apply directly to \mathbf{X} in the multivariate setting. Hence we do not discuss the topic here. Muller and Fetterman (2002, especially Chapters 9, 11, and 16), as well as many other authors of univariate texts, provided detailed discussions of polynomial models.

Natural polynomial models seem easy to handle and interpret. Disadvantages include potential numerical problems and a corresponding lack of statistical independence among the columns of $\widehat{\mathbf{B}}_1$ (which contains the estimators assumed to not be zero, as in the last equation). In turn, dependencies among terms can greatly complicate interpretations (a fact often overlooked). Muller and Fetterman (2002, Chapter 9) addressed the difficulty in the context of comparing added-last and added-in-order tests and corresponding correlations. Orthogonal or orthonormal polynomial models avoid the disadvantages. Within-subject design matrix \mathbf{T} ($m \times p$) is assumed to be rank $m \leq p$. It is orthogonal (by rows) if $\mathbf{T}\mathbf{T}'$ is diagonal and orthonormal (by rows) if $\mathbf{T}\mathbf{T}' = \mathbf{I}_m$.

Definition 13.2 In discussing growth curve models, \mathbf{T}_{NAT} indicates a matrix with $t_{ij} = (d_j)^{i-1}$, a natural polynomial. Similarly, \mathbf{T}_{ORT} indicates $\mathbf{T}_{\text{ORT}}\mathbf{T}'_{\text{ORT}} = \mathbf{D}$, a diagonal matrix, which implies \mathbf{T}_{ORT} is orthogonal by rows. Also, \mathbf{T}_{ORN} indicates a matrix orthonormal by rows, $\mathbf{T}_{\text{ORN}}\mathbf{T}'_{\text{ORN}} = \mathbf{I}$.

Example 13.1 The weight of each of $N = 50$ children is measured in kilograms at 5, 6, and 7 years of age. The resulting $p = 3$ ordered responses for child i is $\mathbf{Y}_i = [\text{WT}_{i5} \text{ WT}_{i6} \text{ WT}_{i7}]$ and $\mathbf{d}' = [5 \ 6 \ 7] = \{d_j\}$. The cohort is assumed to be a single homogeneous group, and consequently the between-individual design matrix has one column, $\mathbf{X} = \mathbf{1}_N = [1 \ \cdots \ 1]'$. Without loss of generality, any \mathbf{X} matrix could be used, but to keep the example simple, we assume no additional predictor variables are needed. We further suppose the data were collected with the understanding the growth curve for the children would be well approximated by a low-order polynomial. Therefore a polynomial matrix \mathbf{T} is desired and may be defined in terms of a natural polynomial:

$$\mathbf{B}_{\text{NAT}} = [\beta_0 \ \beta_1 \ \beta_2] \tag{13.7}$$

and

$$\mathbf{T}_{\text{NAT}} = \begin{bmatrix} 1 & 1 & 1 \\ 5 & 6 & 7 \\ 25 & 36 & 49 \end{bmatrix} = \{d_j^{i-1}\}. \tag{13.8}$$

Element (i, j) of $\mathbf{T} = [\mathbf{t}_1 \ \cdots \ \mathbf{t}_p]$ is a monomial. In turn, the elements of the parameter matrix \mathbf{B} are the coefficients of the natural polynomial. The expected weight of child i at time j is

$$E(y_{ij}|\mathbf{X}, \mathbf{T}) = \mathbf{1}\mathbf{B}_{\text{NAT}}\mathbf{t}_j = [\beta_0 \ \beta_1 \ \beta_2] \begin{bmatrix} 1 \\ d_j \\ d_j^2 \end{bmatrix}. \tag{13.9}$$

If we assume $\beta_2 = 0$, then B is partitioned such that $B_1 = [\beta_0 \ \beta_1]$ and $B_2 = [\beta_2]$. Correspondingly, T is partitioned with $T_1 = \begin{bmatrix} 1 & 1 & 1 \\ 5 & 6 & 7 \end{bmatrix}$ and $T_2 = [25 \ 36 \ 49]$.

Definition 13.3 One-to-one linear transformations exist among natural, orthogonal, and orthonormal polynomial design matrices. A set of matrices, all $p \times p$, may be defined as follows:

$$A_{\text{ORTNAT}} = T_{\text{ORT}} T_{\text{NAT}}^{-1} \tag{13.10}$$

$$A_{\text{ORNNAT}} = T_{\text{ORN}} T_{\text{NAT}}^{-1} \tag{13.11}$$

$$A_{\text{ORNORT}} = T_{\text{ORN}} T_{\text{ORT}}^{-1} \tag{13.12}$$

Although $(T_{\text{NAT}})^{-1}$ always exists, the ratio of its eigenvalues, λ_1/λ_p , approaches zero as $p \rightarrow \infty$. Numerical problems will be encountered when trying to invert large natural polynomial matrices represented in computer arithmetic with typical finite precision. By contrast, inverting the orthonormal polynomial design matrix is perfectly stable and a trivial operation: $(T_{\text{ORN}})^{-1} = (T_{\text{ORN}})'$.

Lemma 13.1 One-to-one linear transformations among natural, orthogonal, and orthonormal polynomial regression coefficients exist. The following one-to-one relationships hold (with all B of dimension $q \times p$ and all A $p \times p$):

$$B_{\text{NAT}} = B_{\text{ORT}} A_{\text{ORTNAT}} \tag{13.13}$$

$$B_{\text{NAT}} = B_{\text{ORN}} A_{\text{ORNNAT}} \tag{13.14}$$

$$B_{\text{ORN}} = B_{\text{ORT}} A_{\text{ORNORT}} \tag{13.15}$$

Proof. Left as an exercise.

Definition 13.4 If the elements of d are equally spaced, then $d_{j+1} - d_j$ is a constant $\forall i$, and an *orthogonal design matrix* T_{ORT} exists with every element an integer. The elements of T_{ORT} can be found in tables of orthogonal polynomial coefficients (in the Appendix, Section A.3).

Example 13.2 For $p = 3$

$$T_{\text{ORT}} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{bmatrix} \begin{array}{l} \leftarrow \text{constant} \\ \leftarrow \text{linear} \\ \leftarrow \text{quadratic} \end{array} \tag{13.16}$$

For $p = 4$

$$\mathbf{T}_{\text{ORT}} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -3 & -1 & 1 & 3 \\ 1 & -1 & -1 & 1 \\ -1 & 3 & -3 & 1 \end{bmatrix} \begin{array}{l} \leftarrow \text{constant} \\ \leftarrow \text{linear} \\ \leftarrow \text{quadratic} \\ \leftarrow \text{cubic} \end{array} \quad (13.17)$$

Dividing t_{ij} by $(\sum_{j=1}^p t^2)^{1/2}$ creates the corresponding orthonormal matrix. Premultiplying by $\text{Dg}[(\sum_{j=1}^p t_{ij}^2)^{-1/2}]$ accomplishes the task.

Example 13.3 For the example,

$$\mathbf{T}_{\text{ORN}} = \begin{bmatrix} 0.577350 & 0.577350 & 0.577350 \\ -0.707107 & 0.000000 & 0.707107 \\ 0.408248 & -0.816497 & 0.408248 \end{bmatrix} \begin{array}{l} \leftarrow \text{constant} \\ \leftarrow \text{linear} \\ \leftarrow \text{quadratic} \end{array} \quad (13.18)$$

In the example of children's weights ($N = 50, m = 2, p = 3, q = 1$), the children's weights were measured at equally spaced time intervals. Therefore an orthogonal matrix exists which has integer entries:

$$\mathbf{T}_{\text{ORT}} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{bmatrix}. \quad (13.19)$$

In turn, the expected value of the weight of child i at time j is

$$\begin{aligned} E(y_{ij} | \mathbf{X} = \mathbf{1}, \mathbf{T}) &= \mathbf{B}_{\text{ORT}} \mathbf{T}_{j\text{ORT}} = [\beta_0 \ \beta_1 \ \beta_2] \begin{bmatrix} -6 + d_j \\ 106 - 36d_j + 3d_j^2 \end{bmatrix} \\ &= [\beta_0 \ \beta_1 \ \beta_2] \begin{bmatrix} 1 & 0 & 0 \\ -6 & 1 & 0 \\ 106 & -36 & 3 \end{bmatrix} \begin{bmatrix} 1 \\ d_j \\ d_j^2 \end{bmatrix}. \end{aligned} \quad (13.20)$$

The \mathbf{B} matrices and \mathbf{T} matrices of the natural and orthogonal models are linearly related. The linear transformation matrix for the two designs is

$$\mathbf{A}_{\text{ORTNAT}} = \begin{bmatrix} 1 & 0 & 0 \\ -6 & 1 & 0 \\ 106 & -36 & 3 \end{bmatrix} = \mathbf{T}_{\text{ORT}} \mathbf{T}_{\text{NAT}}^{-1}. \quad (13.21)$$

Thus we have $\mathbf{B}_{\text{NAT}} = \mathbf{B}_{\text{ORT}} \mathbf{A}_{\text{ORTNAT}}$,

$$[\beta_0 \ \beta_1 \ \beta_2]_{\text{NAT}} = [\beta_0 \ \beta_1 \ \beta_2]_{\text{ORT}} \begin{bmatrix} 1 & 0 & 0 \\ -6 & 1 & 0 \\ 106 & -36 & 3 \end{bmatrix}, \quad (13.22)$$

and $\mathbf{T}_{\text{ORT}} = \mathbf{A}_{\text{ORTNAT}} \mathbf{T}_{\text{NAT}}$,

$$\begin{bmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{bmatrix} = \begin{bmatrix} 1 & 0 & 0 \\ -6 & 1 & 0 \\ 106 & -36 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 \\ 5 & 6 & 7 \\ 25 & 36 & 49 \end{bmatrix}. \tag{13.23}$$

We suppose mean weight is approximately linear in time over the ages from 5 to 7 years. If so, it is appropriate to assume $\beta_2 = 0$ in the natural polynomial model, or $\beta_2 = 0$ in the orthogonal polynomial model. The resulting model in terms of orthogonal polynomials is $E(Y|X, T) = [\beta_0 \ \beta_1]_{\text{ORT}} T_{1,\text{ORT}}$, in which

$$T_{1,\text{ORT}} = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \end{bmatrix}. \tag{13.24}$$

The possibility of specifying a GCM using one of several kinds of polynomials raises the question of equivalence among models more generally in terms of any specifications of the within-individual design matrix T .

Definition 13.5 Two models specified by their design matrices and parameter spaces, $\text{GCM}(Y; X_1 B_1 T_1, \Sigma)$, $B_1 \in \Omega_1$, and $\text{GCM}(Y; X_2 B_2 T_2, \Sigma)$, $B_2 \in \Omega_2$, are *linearly equivalent* \Leftrightarrow
(a) $\forall B_1 \in \Omega_1$ there exists $B_2 \in \Omega_2$ such that $X_1 B_1 T_1 = X_2 B_2 T_2$ and
(b) $\forall B_2 \in \Omega_2$ there exists $B_1 \in \Omega_1$ such that $X_1 B_1 T_1 = X_2 B_2 T_2$.

Having fitted a GCM, it is desirable to create graphical representations of the estimated growth curve. As for all polynomial models, computations are best conducted in orthogonal or orthonormal terms, although displayed in terms of natural polynomials.

Definition 13.6 The *estimated polynomial growth curve*, or *dose-response curve*, is the estimator of mean response as a function of time, or dose, d of the form $\hat{\mu}(d; C) = C \hat{B} u_d$ ($1 \times q \times m \times 1$) with $u_d = [1 \ d \ \dots \ d^{m-1}]'$ ($m \times 1$). A set of coordinates $\{[d, \hat{\mu}(d; C)] : d \in [d_1, d_p]\}$ define an estimated growth curve.

Theorem 13.3 For $\text{GCM}_{N,p,q,m}(Y_i; X_i B_{\text{ORT}} T_{\text{ORT}}, \Sigma)$ and any value of d in the range of the data, a growth curve can be conveniently and plausibly estimated. With C ($1 \times q$) and u_d ($m \times 1$) such that $\{[d, \hat{\mu}(d; C)] : d \in [d_1, d_p]\}$ gives

$$\hat{\mu}(d; C) = C \hat{B}_{\text{ORT}} A_{\text{ORTNAT}} u_d \tag{13.25}$$

with $m \times 1$

$$u_d = [1 \ d \ \dots \ d^{m-1}]'. \tag{13.26}$$

The result allows plotting growth curves (or dose-response curves) at arbitrarily many points in terms of the orthogonal polynomial regression parameters.

Proof. The natural polynomial model and the orthogonal polynomial model are linearly equivalent (with all corresponding matrices of the same dimensions). In particular, $GCM(\mathbf{Y}; \mathbf{X}\mathbf{B}_{\text{ORT}}\mathbf{T}_{\text{ORT}}, \Sigma)$ with $\mathbf{B}_{2\text{ORT}} = \mathbf{0}$ is equivalent to $GCM(\mathbf{Y}; \mathbf{X}\mathbf{B}_{\text{NAT}}\mathbf{T}_{\text{NAT}}, \Sigma)$ with $\mathbf{B}_{2\text{NAT}} = \mathbf{0}$. The linear relationship between natural and orthogonal polynomial coefficients gives $\widehat{\mathbf{B}}_{\text{NAT}} = \widehat{\mathbf{B}}_{\text{ORT}}\mathbf{A}_{\text{ORTNAT}}$. In turn $\widehat{\mu}(d; \mathbf{C}) = \mathbf{C}\widehat{\mathbf{B}}_{\text{NAT}}\mathbf{u}_d = \mathbf{C}(\widehat{\mathbf{B}}_{\text{ORT}}\mathbf{A}_{\text{ORTNAT}})\mathbf{u}_d$. \square

Extrapolating a growth curve, or any other model, outside the range of data can be quite misleading. Muller and Fetterman (2002, Chapter 9) provided examples. As with most issues of practical data analysis and interpretation, we leave the topic for consideration in other settings.

Example 13.4 The integer-valued orthogonal polynomial model for $d \in \{1, 2, 3, 4\}$ is

$$\begin{aligned} E(\mathbf{Y}|\mathbf{X}, \mathbf{T}) &= \mathbf{X}\mathbf{B}_{\text{ORT}}\mathbf{T}_{\text{ORT}} \\ &= \mathbf{X}[\beta_0 \ \beta_1 \ \beta_2 \ \beta_3] \begin{bmatrix} 1 & 1 & 1 & 1 \\ -3 & -1 & 1 & 3 \\ 1 & -1 & -1 & 1 \\ -1 & 3 & -3 & 1 \end{bmatrix}. \end{aligned} \tag{13.27}$$

If we require $\beta_1 = 0$ and $\beta_3 = 0$, then $\mathbf{B}_{\text{ORT}} = [\beta_0 \ \mathbf{0} \ \beta_2 \ 0]$. The corresponding natural polynomial regression coefficients are $\mathbf{B}_{\text{NAT}} = [\beta_0 \ \beta_1 \ \beta_2 \ 0]$. The result is obtained from $\mathbf{B}_{\text{NAT}} = \mathbf{B}_{\text{ORT}}\mathbf{A}_{\text{ORTNAT}}$, with

$$\mathbf{A}_{\text{ORTNAT}} = \begin{bmatrix} 1 & 0 & 0 & 0 \\ -5 & 2 & 0 & 0 \\ 5 & -5 & 1 & 0 \\ -35 & 55.66 & -25 & 3.33 \end{bmatrix}. \tag{13.28}$$

The example natural polynomial has $m = 3$ nonzero coefficients, $\mathbf{B}_{1,\text{NAT}} = [\beta_0 \ \beta_1 \ \beta_2]$, while the orthogonal polynomial has only $m_1 = 2$ nonzero coefficients, $\mathbf{B}_{1,\text{ORT}} = [\beta_0 \ \beta_2]$. The quadratic term of the integer-valued orthogonal polynomial is $\beta_2(5 - d + d^2)$, which includes a linear component. It is important to remember that different kinds of polynomials (of the same order) for a given model may have different numbers of nonzero regression coefficients.

The notation $\mathbf{A}[\cdot, \cdot]$ indicates a submatrix of \mathbf{A} created by selecting all elements of rows in a list and all elements of columns in a list. We use \mathbf{J}_1 ($1 \times m$) to indicate which of p columns are assumed to be *not* zero. With $\mathbf{J}_1 = [1 \ 3]$ and $m_1 = \max\{\mathbf{J}_1\} = 3$, $\mathbf{A}_{\text{ORTNAT}}[\mathbf{J}_1, [1 \ 2 \ \cdots \ m_1]]$ has m rows and m_1 columns and $\widehat{\mathbf{B}}_{\text{INAT}} = \widehat{\mathbf{B}}_{\text{IORT}}\mathbf{A}_{\text{ORTNAT}}[\mathbf{J}_1, [1 \ 2 \ \cdots \ m_1]]$. The \mathbf{J}_1 matrix selects the rows and columns of $\mathbf{A}_{\text{ORTNAT}}$ needed to transform orthogonal polynomial results into results in terms of natural polynomials. Similarly $\mathbf{J}_2 = [2 \ 4]$ ($1 \times p - m$) indicates columns 2 and 4 of the coefficient matrix *are* assumed to be zero.

Theorem 13.4 For $\text{GCM}_{N,p,q,m}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B} \mathbf{T}, \Sigma)$ with \mathbf{B} and \mathbf{T} in terms of an orthogonal or orthonormal polynomial, growth curves can be plotted in terms of \mathbf{J}_1 ($1 \times m$), the list of columns of \mathbf{B} assumed *not* to be zero. At time d the estimated value is specified by \mathbf{C} ($1 \times q$), $\widehat{\mathbf{B}}_1$ ($q \times m$), and \mathbf{A} ($m \times m_1$) as

$$\widehat{\boldsymbol{\mu}}(d; \mathbf{C}) = \mathbf{C} \widehat{\mathbf{B}}_1 \mathbf{A} \mathbf{u}_d \tag{13.29}$$

with $\mathbf{u}_d = [1 \ d \ d^2 \ \cdots \ d^{m-1}]'$ and

$$\mathbf{A} = \mathbf{A}_{\text{ORTNAT}}[\mathbf{J}_1, [1 \ 2 \ \cdots \ m_1]], \tag{13.30}$$

in which $m_1 = \max\{\mathbf{J}_1\}$.

Proof. Left as an exercise.

Corollary 13.4.1 A set of $s > p$ points on a growth curve can be computed with a natural polynomial design matrix with more columns than rows, \mathbf{U} ($p \times s$), the first p rows of \mathbf{T}_{NAT} ($s \times s$). For \mathbf{C} ($1 \times q$), \mathbf{B} ($q \times p$), $\boldsymbol{\mu}'$ ($1 \times s$), \mathbf{d} ($s \times 1$),

$$\widehat{\boldsymbol{\mu}}(d; \mathbf{C}) = \mathbf{C} \widehat{\mathbf{B}}_{\text{ORN}} \mathbf{A}_{\text{ORNNAT}} \mathbf{U}. \tag{13.31}$$

The set cannot be computed with an orthogonal design matrix with $s > p$ columns (the dimension of the fitted model). In particular, $\widehat{\boldsymbol{\mu}}(d; \mathbf{C}) \neq \mathbf{C} \widehat{\mathbf{B}}_{\text{ORN}} \mathbf{U}_{\mathbf{T}}$ if $\mathbf{U}_{\mathbf{T}}$ ($p \times s$) is the first p rows of \mathbf{T}_{ORN} ($s \times s$).

Proof. Left as an exercise.

Corollary 13.4.2 Under the same assumptions, estimated growth curve values can be computed as

$$\widehat{\boldsymbol{\mu}}(d; \mathbf{C}) = \mathbf{C} \widehat{\mathbf{B}}_{\text{ORN}} \mathbf{A}_{\text{ORNNAT}} \mathbf{u}_d, \tag{13.32}$$

with $\mathbf{u}_d = [1 \ d \ d^2 \ \cdots \ d^{m-1}]'$, in which d is any value.

Proof. Left as an exercise.

The example of children's weights ($N = 50$, $q = 1$, $p = 3$, $m = 2$) allows illustrating the process. Having estimated $\mathbf{B}_{\text{ORT}} = [\beta_0 \ \beta_1]$, we can plot the estimated growth curve. Suppose we only want to plot three points on the curve corresponding to $\mathbf{t}' = [d_1 \ d_2 \ d_3]$. With $\mathbf{C} = 1$ and \mathbf{T}_{ORT} (3×3) we have

$$\begin{aligned} \widehat{\boldsymbol{\mu}}(d; \mathbf{C}) &= \mathbf{C} \widehat{\mathbf{B}}_{\text{ORT}} \mathbf{T}_{\text{ORT}} \\ &= 1 [\widehat{\beta}_0 \ \widehat{\beta}_1 \ 0] \mathbf{T}_{\text{ORT}}. \end{aligned} \tag{13.33}$$

Here $\mathbf{C} = 1$ because the participants form a single homogeneous group. The three points correspond to the observed ages (5, 6, and 7 years). Corollary 13.4.1 answers the question "How would additional intermediate points be plotted?"

Example 13.5 We seek to compute five points within the range of the data, corresponding to ages $[5 \ 5.5 \ 6.0 \ 6.5 \ 7.0]' = \mathbf{t}_1$. With $\mathbf{C} = \mathbf{I}$, \mathbf{A} (3×3), \mathbf{U} (3×5), the points would be computed as

$$\begin{aligned}
 \hat{\mu}(\mathbf{t}_1; \mathbf{C}) &= \mathbf{C} \hat{\mathbf{B}}_{\text{ORT}} \mathbf{A}_{\text{ORTNAT}} \mathbf{U} \\
 &= [\hat{\beta}_0 \ \hat{\beta}_1 \ 0] \begin{bmatrix} 1 & 0 & 0 \\ -6 & 1 & 0 \\ 106 & -36 & 3 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 5 & 5.50 & 6 & 6.50 & 7 \\ 25 & 30.25 & 36 & 42.25 & 49 \end{bmatrix} \\
 &= [\hat{\beta}_0 \ \hat{\beta}_1] \begin{bmatrix} 1 & 0 \\ -6 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 5 & 5.5 & 6 & 6.5 & 7 \end{bmatrix} \\
 &= [\hat{\beta}_0 \ \hat{\beta}_1] \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ -1 & -0.5 & 0 & -0.5 & 1 \end{bmatrix}. \tag{13.34}
 \end{aligned}$$

Matrix \mathbf{U} can be interpreted as being \mathbf{T}_{NAT} (3×3) augmented with two additional columns. It can also be interpreted as being the first three rows of \mathbf{T}_{NAT} ($s \times s$). Arbitrarily large numbers of points on the curve may be plotted by augmenting the \mathbf{T}_{NAT} matrix with additional columns. It is usually desirable (but not necessary) to choose evenly spaced points.

13.5 WITHIN-INDIVIDUAL DESIGN MATRICES

Two of the most desirable attributes of the \mathbf{T} matrix are full rank and orthogonality. The creative analyst can draw from the great variety of orthogonal designs to formulate an appropriate GCM. Any full-rank nonorthogonal design matrix can be reparameterized as an orthogonal design. After applying the new design, hypotheses about the original nonorthogonal design are easily tested.

The process may be described in terms of \mathbf{T} , a $p \times p$ arbitrary full-rank design matrix. Since $\mathbf{T}\mathbf{T}'$ is full rank, $p \times p$, and symmetric, a lower triangular matrix \mathbf{L} ($p \times p$) exists such that $\mathbf{T}\mathbf{T}' = \mathbf{L}\mathbf{L}'$ (the Cholesky decomposition). Lower triangular \mathbf{L}^{-1} maps \mathbf{T} onto an orthogonal matrix $\mathbf{T}_{\text{new}} = \mathbf{L}^{-1}\mathbf{T}$. If a model which assumes $E(\mathbf{Y}|\mathbf{X}, \mathbf{T}_{\text{new}}) = \mathbf{X}\mathbf{B}_{\text{new}}\mathbf{T}_{\text{new}}$ is fitted, the original hypothesis $H_0 : \mathbf{C}\mathbf{B}\mathbf{U} = \mathbf{0}$ ($a \times b$) can be tested in the form $H_0 : \mathbf{C}\mathbf{B}_{\text{new}}\mathbf{U}_{\text{new}} = \mathbf{0}$ ($a \times b$), in which $\mathbf{U}_{\text{new}} = \mathbf{L}^{-1}\mathbf{U}$. If some columns of \mathbf{B} ($q \times p$) are assumed to be zero, then \mathbf{T} ($p \times p$) is replaced by \mathbf{T}_1 ($m \times p$) in the above equations.

Polynomials are useful when little is known about the mechanisms underlying the ordered responses. If more is known about the underlying process, then a formulation which uses the additional information likely will be better. As an example, if the underlying process is known to be periodic, then the design matrix \mathbf{T} might be formulated in terms of sine and cosine functions of time.

Definition 13.7 It may be convenient to define transformation matrix T in a GCM in terms of a *discrete Fourier transform* (DFT). If p is even, then $p_c = (p - 0)/2$ and $p_b = (p - 2)/2$. If p is odd, then $p_a = p_b = (p - 1)/2$. For time d_j , $c_{ij} = \cos(2\pi i d_j f)$ and $s_{ij} = \sin(2\pi i t_j f)$, with f the *fundamental frequency* of the process and $1/f$ the *fundamental period*. In turn,

$$T_{\text{DFT}} = \begin{bmatrix} \mathbf{1}' \\ \mathbf{T}_c \\ \mathbf{T}_s \end{bmatrix} \begin{array}{ll} \leftarrow \text{constants} & (1 \times p) \\ \leftarrow \text{cosines} & (p_c \times p) \\ \leftarrow \text{sines} & (p_b \times p) \end{array} \quad (13.35)$$

in which $T_c = \{c_{ij}\}$ with $1 \leq i \leq p_a$ and $T_s = \{s_{ij}\}$ with $1 \leq i \leq p_b$. Larger values of i indicate the higher order terms. The terms with $i = 1$ represent the lowest frequency component. The terms with $i = 2$ represent the *first harmonic*, while $i = 3$ terms give the *second harmonic*.

Theorem 13.5 If $d' = [1 \ 2 \ \cdots \ p] = \{d_j\} = \{j\}$ and $f = 1/p$, then the DFT design matrix T_{DFT} is a $p \times p$ orthogonal matrix.

Proof. Left as an exercise.

So far we have assumed the measurement made at time j is a function of d_j , which implies the ordered columns of Y ($N \times p$) are functionally related. The design matrix T can, when necessary, be formulated in terms of a polynomial in more than one variable.

Example 13.6 Temporal changes in level of pain (Y) indexed by time (t) in hours can be observed for several doses (d) of an analgesic drug via a multi-period crossover design; i.e., the individual sequentially experiences the different doses in different periods of the study and longitudinal observations are recorded within each period. In other human studies, simultaneous administration of different doses of a drug can occur, for example, when two different doses of a drug are applied to the left and right eyes respectively. Other examples include allergy skin testing (various patches of skin are exposed to different allergens) and in vitro studies [e.g., a blood sample is drawn, divided into multiple subsamples of equal volume (aliquots), each aliquot is treated with one of the dose levels of interest, and response in each aliquot is then observed longitudinally].

For illustration, suppose two within-subject factors, time (t) and dose (d), are present, and further suppose a response y is observed at $p = 3$ times for each of 3 doses. In terms of a two-variable natural polynomial the mean response for the 3×3 factorial experiment is of the form

$$\begin{aligned} \mu_{ijk} &= E[y_i(t_j, d_k) | \mathbf{X}, \mathbf{T}]_j \\ &= \mathbf{X}_i \mathbf{B} t_j, \end{aligned} \quad (13.36)$$

in which \mathbf{T} is the Kronecker product of two single-variable design matrices:

$$\mathbf{T}(t, d) = \mathbf{T}(t) \otimes \mathbf{T}(d). \tag{13.37}$$

If the participants are a homogeneous group, then $\mathbf{X} = \mathbf{1}$ ($N \times 1$). If the two-variable polynomial is of maximal order, then

$$\mu_{ijk} = \beta_{00} + \beta_{10}t + \beta_{20}t^2 + \beta_{01}d + \beta_{11}td + \beta_{21}t^2d + \beta_{02}d^2 + \beta_{12}td^2 + \beta_{22}t^2d^2. \tag{13.38}$$

For $\mathbf{T}(t, d)$ to be an orthogonal matrix, a Kronecker product of orthogonal \mathbf{T} matrices is used: $\mathbf{T}(t, d) = \mathbf{T}(t) \otimes \mathbf{T}(d)$ is (9×9) and

$$\mathbf{T}(t) = \mathbf{T}(d) = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 0 & 1 \\ 1 & -2 & 1 \end{bmatrix}. \tag{13.39}$$

Similarly the design matrix for a polynomial in three variables is

$$\mathbf{T}(t, s, u) = \mathbf{T}(t) \otimes \mathbf{T}(s) \otimes \mathbf{T}(u). \tag{13.40}$$

In general, the \mathbf{T} matrix for a polynomial in K variables is the Kronecker product of K matrices. The inverse of \mathbf{T} is the Kronecker product of the individual inverses. If $\mathbf{T}(t, d) = \mathbf{T}(t) \otimes \mathbf{T}(d)$, then $\mathbf{T}^{-1}(t, d) = \mathbf{T}^{-1}(t) \otimes \mathbf{T}^{-1}(d)$.

In the GCM, the definition of matrix \mathbf{X} ($N \times q$) does not require any additional considerations beyond the ones made in the multivariate GLM. Therefore the between-individual design matrix \mathbf{X} may be thought of as the usual \mathbf{X} matrix in a general linear model. Given a matrix of responses \mathbf{Y} ($N \times p$), the definition of matrix \mathbf{T} ($p \times p$) does not place any restrictions on the definition of \mathbf{X} .

13.6 ESTIMATION METHODS

Under the assumptions of $\text{GCM}_{N,p,q,m}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}\mathbf{T}, \Sigma)$, the correct model for a single individual's responses is

$$\text{E}(\mathbf{Y}_i' | \mathbf{X}_i, \mathbf{T}) = (\boldsymbol{\theta}_i' \mathbf{T})', \tag{13.41}$$

with $m \times 1$ $\boldsymbol{\theta}_i = (\mathbf{X}_i\mathbf{B})'$ and $\mathcal{V}(\mathbf{Y}_i' | \mathbf{X}_i, \mathbf{T}) = \Sigma$ ($p \times p$). Multiplying both sides of the regression equation yields

$$\text{E}[\mathbf{Y}_i\mathbf{V}^{-1}\mathbf{T}'(\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1} | \mathbf{X}_i, \mathbf{T}] = \boldsymbol{\theta}_i' = \mathbf{X}_i\mathbf{B}. \tag{13.42}$$

The vector $\hat{\boldsymbol{\theta}}_i = (\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}\mathbf{T}\mathbf{V}^{-1}\mathbf{Y}_i'$ is an $m \times 1$ weighted-least-squares estimator (unweighted if $\mathbf{V} = \mathbf{I}_p$) for the individual-specific model $\text{E}(\mathbf{Y}_i' | \mathbf{X}_i, \mathbf{T}) = \mathbf{T}'\boldsymbol{\theta}_i$. Given $\boldsymbol{\theta}_i' = \mathbf{X}_i\mathbf{B}$, with $\mathbf{X}_i\mathbf{B}$ invariant to the choice of \mathbf{V} , it follows $\boldsymbol{\theta}_i$ is invariant to the choice of \mathbf{V} . Also

$$\begin{aligned} E(\hat{\theta}_i) &= E[(\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}\mathbf{T}\mathbf{V}^{-1}\mathbf{Y}'_i] \\ &= (\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}\mathbf{T}\mathbf{V}^{-1}E(\mathbf{Y}'_i) \\ &= (\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}\mathbf{T}\mathbf{V}^{-1}(\theta'_i\mathbf{T})' \\ &= \theta_i. \end{aligned} \tag{13.43}$$

Although $\hat{\theta}'_i$ is unbiased for all choices of \mathbf{V} , its precision is optimal if $\mathbf{V} = \Sigma$. Unfortunately Σ is unknown.

If a separate univariate GLM model is fitted for each individual, it would then be reasonable to fit a multivariate GLM for $\hat{\Theta} = [\hat{\theta}_1 \cdots \hat{\theta}_N]' = \mathbf{Y}\mathbf{V}^{-1}\mathbf{T}(\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}$ ($N \times m$), with $N \times q \times m$ mean

$$E(\hat{\theta}|\mathbf{X}) = \mathbf{X}\mathbf{B} \tag{13.44}$$

and $\mathcal{V}[\text{vec}(\hat{\theta})|\mathbf{X}] = \Omega \otimes \mathbf{I}$. Here $\Omega = \mathbf{G}'\Sigma\mathbf{G}$ with $\mathbf{G} = \mathbf{V}^{-1}\mathbf{T}(\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}$ gives estimator

$$\hat{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\hat{\Theta} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}\mathbf{V}^{-1}\mathbf{T}(\mathbf{T}\mathbf{V}^{-1}\mathbf{T}')^{-1}. \tag{13.45}$$

The variance of the estimator,

$$\mathcal{V}[\text{vec}(\hat{\mathbf{B}})|\mathbf{X}] = (\mathbf{G}'\Sigma\mathbf{G}) \otimes (\mathbf{X}'\mathbf{X})^{-1}, \tag{13.46}$$

depends on the choice of \mathbf{V} . Obvious choices for \mathbf{V} include $\tilde{\Sigma} = \mathbf{Y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y}/N$, the MLE, and the unbiased estimator $\hat{\Sigma} = \tilde{\Sigma}N/[N - \text{rank}(\mathbf{X})]$.

The results just stated could be derived in terms of a transformation of the columns of an ordinary multivariate GLM using $\mathbf{G} = \mathbf{T}^{-1}$ ($p \times p$), combined with the fact that conditioning on \mathbf{T} is the same as conditioning on \mathbf{G} . In particular,

$$E(\mathbf{Y}\mathbf{G}|\mathbf{X}, \mathbf{G}) = \mathbf{X}\mathbf{B} \tag{13.47}$$

and

$$\mathcal{V}(\mathbf{Y}\mathbf{G}|\mathbf{X}, \mathbf{G}) = (\mathbf{G}'\Sigma\mathbf{G}) \otimes \mathbf{I}_N. \tag{13.48}$$

The transformation yields new dependent variables, $\mathbf{Y}\mathbf{G}$, and provides a GLM setting for estimation and inference. The partitioned model

$$\begin{aligned} E(\mathbf{Y}|\mathbf{X}, \mathbf{T}) &= \mathbf{X}\mathbf{B}\mathbf{T} \\ &= \mathbf{X}[\mathbf{B}_1 \quad \mathbf{B}_2] \begin{bmatrix} \mathbf{T}_1 \\ \mathbf{T}_2 \end{bmatrix} \end{aligned} \tag{13.49}$$

induces the partitioning $\mathbf{Y}\mathbf{G} = [\mathbf{Y}\mathbf{G}_1 \quad \mathbf{Y}\mathbf{G}_2] = [\mathbf{Y}_1 \quad \mathbf{Y}_2]$. Constrained by $\mathbf{B}_2 = \mathbf{0}$, here $E(\mathbf{Y}_2|\mathbf{X}, \mathbf{T}) = \mathbf{0}$ [$N \times (p - q)$] while $E(\mathbf{Y}_1|\mathbf{X}, \mathbf{T}) = \mathbf{X}\mathbf{B}_1$. The columns of \mathbf{G}_1 are usually the first few columns of \mathbf{G} . In any case, the rows of \mathbf{T} can always be permuted (along with columns of \mathbf{B} and columns of \mathbf{G}) so no generality is lost by assuming \mathbf{G}_1 contains the first few columns of \mathbf{G} ($p \times p$).

Restricting some columns of \mathbf{B} to be zero is achieved by omitting \mathbf{Y}_2 ($N \times p - q$) from the model. The test of $H_0 : \mathbf{B}_2 = \mathbf{0}$ serves as a goodness-of-fit test.

If all of the \mathbf{Y}_2 variables are omitted from the GLM in order to restrict \mathbf{B} ($q \times p$), then the following model is fitted:

$$E(\mathbf{Y}_1 | \mathbf{X}, \mathbf{G}_1) = \mathbf{X}\mathbf{B}_1 \tag{13.50}$$

$$\mathcal{V}(\mathbf{Y}_1 | \mathbf{X}, \mathbf{G}_1) = \mathbf{G}'_1 \boldsymbol{\Sigma} \mathbf{G}_1 \otimes \mathbf{I}. \tag{13.51}$$

Rao (1965) and Khatri (1966) described how to reduce a GCM to an ordinary multivariate GLM. They proved some information is discarded if the \mathbf{Y}_2 variables are ignored. Furthermore the use of some or all of \mathbf{Y}_2 as independent covariates can improve the power of tests and reduce the widths of confidence intervals. Including none of the \mathbf{Y}_2 variables as covariates yields the unweighted estimator

$$\tilde{\mathbf{B}}_1 = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}\mathbf{T}'_1(\mathbf{T}_1\mathbf{T}'_1)^{-1}, \tag{13.52}$$

while using all of the \mathbf{Y}_2 variables as covariates yields the weighted estimator

$$\hat{\mathbf{B}}_1 = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}\hat{\boldsymbol{\Sigma}}^{-1} \mathbf{T}'_1(\mathbf{T}_1\hat{\boldsymbol{\Sigma}}^{-1} \mathbf{T}'_1)^{-1}, \tag{13.53}$$

in which $\hat{\boldsymbol{\Sigma}}$ is the MLE of $\boldsymbol{\Sigma}$. Using all \mathbf{Y}_2 variables requires fitting the model

$$E(\mathbf{Y}_1 | \mathbf{X}, \mathbf{G}_1, \mathbf{Y}_2) = \mathbf{X}\mathbf{B}_1 + \mathbf{Y}_2\boldsymbol{\Gamma} \tag{13.54}$$

$$\mathcal{V}(\mathbf{Y}_1 | \mathbf{X}, \mathbf{G}_1, \mathbf{Y}_2) = \mathbf{G}'_1 \boldsymbol{\Sigma} \mathbf{G}_1 \otimes \mathbf{I}, \tag{13.55}$$

which is equivalent to computing $\hat{\mathbf{B}}_1 = (\mathbf{X}'\mathbf{X})^{-1} \mathbf{X}'\mathbf{Y}\mathbf{V}^{-1} \mathbf{T}'(\mathbf{T}\mathbf{V}^{-1} \mathbf{T}')^{-1}$ with $\mathbf{V} = \hat{\boldsymbol{\Sigma}}$. Grizzle and Allen (1969) derived the variance of $\hat{\mathbf{B}}_1$:

$$\mathcal{V}[\text{vec}(\hat{\mathbf{B}}) | \mathbf{X}] = \mathbf{G}'\boldsymbol{\Sigma}\mathbf{G} \otimes (\mathbf{X}'\mathbf{X})^{-1} (N - 1) / [N - (p - q) - 1]. \tag{13.56}$$

Conditioning on some or all of the \mathbf{Y}_2 variables uses additional degrees of freedom. Since some of the \mathbf{Y}_2 variables may be redundant of one another in the information they provide, Rao (1965) and Grizzle and Allen (1969) proposed balancing the gain of information against the use of degrees of freedom by including only a few well-chosen \mathbf{Y}_2 covariates. As criteria for including or excluding variables, Rao proposed examining the resulting widths of confidence intervals for elements of $\hat{\mathbf{B}}$. Grizzle and Allen proposed relying on a measure of generalized variance, the determinant of the covariance matrix of estimator $\hat{\mathbf{B}}$. For confirmatory hypothesis testing, the conservative approach is to include all of the \mathbf{Y}_2 covariates.

Berger (1986) conducted simulations evaluating various strategies for growth curve analysis. His results indicate that using observed properties of the data to help model the covariance inflates the type I error rate. Hence he recommended avoiding the approach unless an appropriate correction could be determined.

13.7 RELATIONSHIPS TO THE UNIVARIATE AND MIXED MODELS

The model $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$ has the corresponding equation for sampling unit i :

$$\mathbf{Y}'_i = (\mathbf{X}_i \otimes \mathbf{I}_p)\text{vec}(\mathbf{B}') + \mathbf{E}'_i. \tag{13.57}$$

For $p \times m$ \mathbf{T} the model $\mathbf{Y} = \mathbf{X}\mathbf{B}\mathbf{T} + \mathbf{E}$ has corresponding equation for sampling unit i :

$$\mathbf{Y}'_i = (\mathbf{X}_i \otimes \mathbf{I}_m)\text{vec}(\mathbf{T}'\mathbf{B}') + \mathbf{E}'_i. \tag{13.58}$$

From Theorem 1.5, $\text{vec}(\mathbf{ABC}) = (\mathbf{C}' \otimes \mathbf{A})\text{vec}(\mathbf{B})$. With $\mathbf{T}'\mathbf{B}' = \mathbf{T}'\mathbf{B}'\mathbf{I}_q$,

$$\mathbf{Y}'_i = (\mathbf{X}_i \otimes \mathbf{I}_m)(\mathbf{I}_q \otimes \mathbf{T}')\text{vec}(\mathbf{B}') + \mathbf{E}'_i. \tag{13.59}$$

Hence

$$\begin{aligned} \mathbf{Y}'_i &= (\mathbf{X}_i \otimes \mathbf{T}')\text{vec}(\mathbf{B}') + \mathbf{E}'_i. & (13.60) \\ (m \times 1) &= [(1 \times q) \otimes (m \times p)](pq \times 1) + (m \times 1) \\ (m \times 1) &= (m \times qp)(pq \times 1) + (m \times 1) \end{aligned}$$

Here \mathbf{X}_i contains between-subject design information, while \mathbf{T} contains within-subject design information. The last equation corresponds directly to a mixed model or a GGLM. In contrast to a GCM, a mixed model allows elements of \mathbf{T} and m to vary across subjects.

It can be extremely helpful in understanding an incomplete design to examine the corresponding complete (factorial design). Considering restrictions or variable deletions can then be used to specify the incomplete design in terms of the complete design. The model form just described is intended to help the reader in such an endeavor in the context of mixed models.

EXERCISES

13.1 An investigator gathered n_i observations on each of N sampling units for $i \in \{1, \dots, N\}$ from two different groups. A total of N_1 sampling units were selected from group 1 and N_2 sampling units were from group 2 ($N_1 + N_2 = N$). The n_i observations from each sampling unit are assumed to be independent of the observations from another sampling unit and follow the linear model $\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\gamma} + \mathbf{e}_{iy}$ with $\mathbf{X}_i = [\mathbf{1}_{n_i} \quad \mathbf{0}]$ if sampling unit i belongs to group 1 and $\mathbf{X}_i = [\mathbf{0} \quad \mathbf{1}_{n_i}]$ if sampling unit i belongs to group 2. Also $\boldsymbol{\gamma} = [\gamma_1 \quad \gamma_2]'$ is fixed and unknown, while $\mathbf{e}_{iy} \sim N_{n_i}(\mathbf{0}, \sigma^2\mathbf{I}_{n_i})$ with σ^2 unknown. The original data, y_{ij} , are not available, but we do know $\sum_{j=1}^{n_i} y_{ij}$ for each sampling unit. From the sums, we calculate $m_i = \sum_{j=1}^{n_i} y_{ij}/n_i$ and consider the following linear model:

$$\mathbf{m} = \{\mathbf{m}_i\} = \mathbf{Z}\boldsymbol{\beta} + \mathbf{e}, \text{ with } \boldsymbol{\beta} = [\beta_1 \quad \beta_2] \text{ and } \mathbf{Z} = \begin{bmatrix} \mathbf{1}_{N_1} & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{N_2} \end{bmatrix}.$$

Hint for some parts of the exercise: Summation notation may help simplify.

13.1.1 Completely specify the distribution of \mathbf{e} .

- 13.1.2 Find an expression for the OLS estimator $\hat{\beta}$. Simplify the expression, and very briefly describe the nature of the elements of $\hat{\beta}$ in terms understandable to a scientific collaborator.
- 13.1.3 Find $E(\hat{\beta})$. Simplify the expression.
- 13.1.4 Find $\mathcal{V}(\hat{\beta})$. Simplify the expression.
- 13.1.5 Enough is known to allow using (exact) weighted least squares. Given what is known, provide an expression for an appropriate WLS estimator $\tilde{\beta}$ and simplify it.
- 13.1.6 Specify the model for m associated with 13.1.5 as a GGLM.
- 13.1.7 Create a linearly equivalent GLM for the GGLM in 13.1.6.
- 13.1.8 Find $E(\tilde{\beta})$. Simplify the expression.
- 13.1.9 Find $\mathcal{V}(\tilde{\beta})$. Simplify the expression.
- 13.1.10 Compare $\mathcal{V}(\hat{\beta})$ with $\mathcal{V}(\tilde{\beta})$. In particular, is $\mathcal{V}(\hat{\beta}) - \mathcal{V}(\tilde{\beta})$ always positive definite, negative definite, nonnegative definite, or nonpositive definite? Explain why briefly. Based on your results in the exercises, can either or both be ruled out as a candidate for being the BLUE?

CHAPTER 14

Estimation for Linear Mixed Models

14.1 MOTIVATION

Given the background of the previous chapters, the theory of estimation may be stated quite simply for the general linear mixed model. Without the need to assume any particular distribution, the method of approximate weighted least squares provides parameter estimates. With the additional assumption of Gaussian errors, iterating the process yields maximum likelihood estimates.

Although the mixed model theory of estimation takes simple forms, successfully finding numerical solutions satisfying the theory sometimes proves difficult with real data. The difficulty often has the unfortunate effect of encouraging practitioners to choose a simple covariance model without data to support the choice. Instead, the analyst should first round the within-subject predictor values to the most coarse resolution scientifically acceptable (often time has been recorded with much numerical precision) before attempting analysis. Second, the analyst should improve the scaling, centering, and coding of the data, as discussed in Chapters 8, 9, and 12 in Muller and Fetterman (2000). If a model with an appropriate and defensible covariance structure will not converge, a different analysis strategy, other than oversimplifying the covariance structure, should be considered. In particular, compound symmetry does not seem to be a plausible model for many types of repeated measures in time.

Throughout most of our discussion of linear models, we make two strong assumptions: a valid model and a sample size sufficient to compute estimates. In contrast to univariate and multivariate model properties, the two assumptions do not suffice to guarantee optimal estimators for mixed models. We must also require symmetry of distributions to guarantee unbiased estimates of expected-value parameters. More importantly, unbiased estimates of covariance parameters are available only with complete and balanced designs (no missing or mistimed data). In practice, such models very often correspond to multivariate models, which can and should be analyzed with multivariate techniques in order to use the best available methods for estimation and inference. Some limitations of the mixed model gradually diminish as sample size increases.

Given that relatively little is known about the finite-sample properties of mixed model methods, we should consider them to be “large-sample” methods. The question naturally arises: How big is “large?” We speculate that serious concerns

about accuracy may be present with fewer than $N = 100$ independent sampling units (ISUs), and no more than a modest number of observations for each (perhaps all $p_i < 10$). It should be emphasized that, except in very special cases, simply having a large number of *observations* ($n = \sum_{i=1}^N p_i$) is not reassuring. The ratio of N to n , as well as the absolute size, plays a role in the performance of mixed models.

14.2 STATEMENT OF THE GENERAL LINEAR MIXED MODEL

The definition and notation of the mixed model was introduced in Chapter 5.

14.3 ESTIMATION AND ESTIMABILITY

Except for various special cases, the mixed model does not have closed-form expressions for estimators for any criterion. Gaussian theory estimation procedures for the mixed model with incomplete and unbalanced data include maximum likelihood (ML), restricted maximum likelihood (REML), moment estimators, and general linear model (ANOVA) estimators. Hocking (1985) discussed estimation for each. Harville (1977) gave a comprehensive review of ML and REML procedures, along with computational techniques. Laird and Ware (1982) discussed the Bayesian approach to variance component estimation, its relationship to REML estimation, and the application of the EM algorithm. Fairclough and Helms (1986) and Andrade and Helms (1986) explored ML estimation for the mixed model with linear covariance structure. More recently, Vonesh and Chinchilli (1997), Verbeke and Molenberghs (2000), and Demidenko (2004) provided book-length treatments of mixed models.

Conditional on knowing $\{\Sigma_i\}$ *exactly*, the theory of exact weighted least squares applies, which implies the theory of estimability and linearly equivalent models developed in previous chapters applies. Our treatment of estimability and linear equivalence concerned only the expected-value parameters $\{\beta_j\} = \beta$ (the fixed effect parameters). The possibility of ambiguous parameter sets caused by (1) purposefully or accidentally LTFR between-subject (\mathbf{X}) design matrices or (2) explicit restrictions on rows of β motivated the developments for univariate models. Multivariate and related growth curve models introduce the additional concerns of (3) purposefully or accidentally LTFR within-subject design matrices (\mathbf{T}, \mathbf{U}) or (4) explicit restrictions on columns of the multivariate form of the expected value parameters, $\mathbf{B} = [\beta_1 \cdots \beta_p]$. The mixed model form may be created by vertically concatenating all of the columns of \mathbf{B} to give $\beta_{\text{mixed}} = \text{vec}(\mathbf{B}) = [\beta_1' \cdots \beta_p']'$. All aspects of the same estimability issues arise in mixed models, albeit with between-subject and within-subject design and parameters jumbled together.

The simplicity of the GLM covariance model essentially eliminates any need to study the analog of estimability for covariance parameters. The (univariate)

GLM_{*N,q*}(*y_i*; **X_i****β**, σ²) describes a model for *N* ISUs and 1 observation per ISU. Stacking all of the data together gives an *N* × 1 vector **y** with **y**' = [*y*'₁ ⋯ *y*'_{*N*}]. The associated covariance model matrix for the entire set of data describes a total of *n* = *N* observations,

$$\begin{aligned} \mathcal{V}(\mathbf{y}) &= \mathbf{I}_N \sigma^2 = \mathbf{I}_N \otimes \sigma^2 \\ &= \begin{bmatrix} \sigma^2 & & 0 \\ & \ddots & \\ 0 & & \sigma^2 \end{bmatrix}. \end{aligned} \tag{14.1}$$

The (multivariate) GLM_{*N,p,q*}(**Y_i**; **X_i****B**, **Σ**) describes a model for *N* ISUs and *p* observations per ISU. Stacking all of the data together by participant gives an (*Np*) × 1 vector, **vec(Y')** with **vec(Y')** = [**Y**'₁ ⋯ **Y**'_{*N*}]. The associated covariance model for all data describes a total of *n* = *Np* observations,

$$\begin{aligned} \mathcal{V}[\mathbf{vec}(\mathbf{Y}')] &= \mathbf{I}_N \otimes \mathbf{\Sigma} \\ &= \begin{bmatrix} \mathbf{\Sigma} & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{\Sigma} \end{bmatrix}. \end{aligned} \tag{14.2}$$

Although rarely considered, the multivariate GLM theory of estimation can tolerate a singular **Σ**, if handled carefully. The calculation of primary parameter estimates (**B̂**, **Σ̂**) does also. An appropriate choice of **U** matrix or multivariate restrictions (**R_rB****R_y** = **A**) avoids any difficulties in testing associated with defining a singular error covariance matrix **Σ_{*}** = **U'ΣU**.

In contrast to the univariate and multivariate GLM, the general linear mixed model LMM_{*N,p_i,q,m*}(**y_i**; **X_i****β**, **Z_i****Σ_{di}**(**τ_d**)**Z_i**' + **Σ_{ei}**(**τ_e**)] describes a model for *N* ISUs and *p_i* observations per ISU. The presence of subscript *i* allows the number of observations to vary across ISU. Stacking all data together by participant gives, with *n* = ∑_{*i*=1}^{*N*} *p_i*, an *n* × 1 vector **y_s** with **y**'_{*s*} = [**y**'₁ ⋯ **y**'_{*N*}]. The associated covariance model for the entire set of data describes a total of *n* observations,

$$\begin{aligned} \mathbf{\Sigma}_s = \mathcal{V}(\mathbf{y}_s) &= \bigoplus_{i=1}^N (\mathbf{Z}_i \mathbf{\Sigma}_{di} \mathbf{Z}_i' + \mathbf{\Sigma}_{ei}) \\ &= \bigoplus_{i=1}^N \mathbf{\Sigma}_i \\ &= \begin{bmatrix} \mathbf{\Sigma}_1 & & \mathbf{0} \\ & \ddots & \\ \mathbf{0} & & \mathbf{\Sigma}_N \end{bmatrix}. \end{aligned} \tag{14.3}$$

Edwards, Stewart, Muller, and Helms (2001) described how linear restrictions can be defined which allow identifying a corresponding unrestricted and linearly equivalent general mixed model. Their results apply to linear constraints providing parallel restrictions on the fixed effect (expected-value) parameters and random effect (covariance) parameters and the class of two-stage mixed models. Hence

much work remains to be done to extend the results to the full range of mixed models and constraints.

Heterogeneity among $\{\Sigma_i\}$ leads to heterogeneity among $\{\widehat{\Sigma}_i\}$, which serve as weights for $\widehat{\beta}$ during the *iterative* calculation of estimates. Despite having full-rank $\{\Sigma_i\}$, either $\{\widehat{\Sigma}_i\}$ or expressions based on the estimates may have difficulty maintaining their intended ranks due to finite precision computer arithmetic. Imprecision may arise due to limitations of the sample size, timing of observations, poor scaling, lack of centering, or poor predictor coding. Estimability may fail at any iteration. Much remains to be learned about how to overcome such local difficulties.

14.4 SOME SPECIAL TYPES OF MODELS

Practical use of the linear mixed model requires further assumptions and constraints to reduce the number of covariance parameters to a manageable level. For stationary processes we often assume that the j, k element of the covariance matrix, $\sigma_{jk} = \langle \Sigma_i(\tau) \rangle_{jk}$, is a simple function of the elapsed times between the various pairs of observations. In particular, an autoregressive covariance model of order 1, AR(1), assumes $\langle \Sigma_i(\tau) \rangle_{jk} = \sigma^2 \rho^{|t_j - t_k|}$ for observation j at time t_j and observation k at time t_k . The AR(1) structure specifies a nonlinear model for the variances and covariances. It has two parameters, $\tau' = [\sigma^2 \rho]$. In some cases, an inherently linear model for the covariances can be defined.

Definition 14.1 A LMM $_{N,p_i,q,m}[\mathbf{y}_i; \mathbf{X}_i\beta, \mathbf{Z}_i\Sigma_{di}(\tau_d)\mathbf{Z}'_i + \Sigma_{ei}(\tau_e)]$ with

$$\begin{aligned} \Sigma_i(\tau) &= \Sigma_{di}(\tau_d)\mathbf{Z}'_i + \Sigma_{ei}(\tau_e) \\ &= \sum_{k=1}^K \tau_k \mathbf{G}_{ik} \end{aligned} \tag{14.4}$$

and all $\{\mathbf{G}_{ik}\}$ known constants which may vary with i has *linear covariance structure*.

A wide variety of useful covariance models can be expressed as a linear structure. Andrade and Helms (1986) developed estimators and test statistics for linear hypotheses on β and τ under the assumption of linear covariance structure.

Example 14.1 If $\mathbf{Z}_i = \mathbf{1}_{p_i}$ while $\Sigma_{di}(\tau_d) = \tau_1$ and $\Sigma_{ei}(\tau_e) = \tau_2 \mathbf{I}_{p_i}$, then $\Sigma_i(\tau) = \tau_1 \mathbf{1}_{p_i} \mathbf{1}'_{p_i} + \tau_2 \mathbf{I}_{p_i}$ exhibits linear covariance structure. The covariance model corresponds to assuming compound symmetry for each independent sampling unit. What are ρ and σ^2 as functions of $\{\tau_1, \tau_2\}$?

Definition 14.2 A LMM $_{N,p_i,q,m}[\mathbf{y}_i; \mathbf{X}_i\beta, \mathbf{Z}_i\Sigma_{di}(\tau_d)\mathbf{Z}'_i + \Sigma_{ei}(\tau_e)]$ with full-rank $\mathbf{X}_i \equiv \mathbf{Z}_i \equiv \mathbf{X}_1$ is a *balanced random coefficient model*.

Demidenko (2004, p. 62) noted that a balanced random coefficient model corresponds exactly to a particular form of a growth curve model. Hence the machinery for the $GLM_{N,p}()$ multivariate model can provide noniterative and optimal estimators as well as exact size- α tests in small samples. Except for (even more) special cases, the commonly used mixed model tests will not coincide with the optimal (multivariate) tests. Although we do not pursue the topic here, we note that a wider class of mixed models may be cast and analyzed as multivariate models and an even wider set as multivariate models with missing data.

14.5 ML ESTIMATION

As discussed in Chapter 5, the mixed model for ISU i may be expressed as $\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i$, with $E(\mathbf{y}_i) = \mathbf{X}_i\boldsymbol{\beta}$ and $\mathcal{V}(\mathbf{y}_i) = \boldsymbol{\Sigma}_i(\boldsymbol{\tau})$ abbreviated as $\boldsymbol{\Sigma}_i$. In turn,

$$\begin{aligned} \mathbf{e}_{+i} &= \mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta} \\ &= \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i. \end{aligned} \tag{14.5}$$

Recalling that $n = \sum_{i=1}^N p_i$, the joint log likelihood is

$$\log L(\boldsymbol{\beta}, \boldsymbol{\tau}) = -\frac{1}{2}n\log(2\pi) - \frac{1}{2}\sum_{i=1}^N [\log|\boldsymbol{\Sigma}_i| + (\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})'\boldsymbol{\Sigma}_i^{-1}(\mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta})] \tag{14.6}$$

or

$$-2\log L(\boldsymbol{\beta}, \boldsymbol{\tau}) = n\log(2\pi) + \sum_{i=1}^N [\log|\boldsymbol{\Sigma}_i| + \text{tr}(\mathbf{e}_{+i}\mathbf{e}_{+i}'\boldsymbol{\Sigma}_i^{-1})]. \tag{14.7}$$

Here $\partial\text{vec}[\boldsymbol{\Sigma}_i]/\partial\boldsymbol{\beta} = \mathbf{0}$. For any fixed $\boldsymbol{\tau}$ the value of $\boldsymbol{\beta}$ which maximizes the likelihood is the weighted least squares estimate. Using $\hat{\boldsymbol{\tau}}$ leads to an approximate weighted least squares estimator,

$$\begin{aligned} \boldsymbol{\beta}(\hat{\boldsymbol{\tau}}) &= (\mathbf{X}'_s\hat{\boldsymbol{\Sigma}}_s^{-1}\mathbf{X}_s)^{-1}(\mathbf{X}'_s\hat{\boldsymbol{\Sigma}}_s^{-1}\mathbf{y}_s) \\ &= \left(\sum_{i=1}^N \mathbf{X}'_i\hat{\boldsymbol{\Sigma}}_i^{-1}\mathbf{X}_i\right)^{-1} \left(\sum_{i=1}^N \mathbf{X}'_i\hat{\boldsymbol{\Sigma}}_i^{-1}\mathbf{y}_i\right). \end{aligned} \tag{14.8}$$

Hence the task reduces to maximizing the profile likelihood, with $\hat{\mathbf{e}}_{+i} = \mathbf{y}_i - \mathbf{X}_i\boldsymbol{\beta}(\hat{\boldsymbol{\tau}})$,

$$\log L_{ML}(\hat{\boldsymbol{\tau}}) = -\frac{1}{2}n\log(2\pi) - \frac{1}{2}\sum_{i=1}^N [\log|\boldsymbol{\Sigma}_i| + \text{tr}(\boldsymbol{\Sigma}_i^{-1}\hat{\mathbf{e}}_{+i}\hat{\mathbf{e}}_{+i}')]. \tag{14.9}$$

Finding the MLE of $\boldsymbol{\tau}$ requires solving $\partial L(\hat{\boldsymbol{\tau}})/\partial\hat{\boldsymbol{\tau}} = \mathbf{0}$. The most useful representations of the system of estimating equations depend on features of the model. If $\boldsymbol{\Sigma}_i$ has linear structure $[\mathcal{V}(\mathbf{y}_i) = \boldsymbol{\Sigma}_i(\boldsymbol{\tau}) = \sum_g \tau_g \mathbf{G}_{ig}$ with $\{\mathbf{G}_{ig}\}$ known],

then the MLE of τ is

$$\hat{\tau} = \left[\left\langle \sum_{i=1}^N \text{tr}(\hat{\Sigma}_i^{-1} \mathbf{G}_{ig} \hat{\Sigma}_i^{-1} \mathbf{G}_{ih}) \right\rangle_{gh} \right]^{-1} \left[\left\langle \sum_{i=1}^N \text{tr}(\hat{\mathbf{e}}_{+i} \hat{\mathbf{e}}'_{+i} \hat{\Sigma}_i^{-1} \mathbf{G}_{ig} \hat{\Sigma}_i^{-1}) \right\rangle_g \right]. \quad (14.10)$$

Here $[\langle \rangle_g]$ is a vector and the expression in $\langle \rangle$ is the value of its element g . Matrix $[\langle \rangle_{gh}]$ is similarly defined and has the indicated value for the element in row g and column h . Advantageous simplifications of the estimating equations can depend on whether the data are balanced, whether the repeated measurements are consistently timed, or whether compound symmetry is assumed.

Except for special cases, the above equations for $\hat{\beta}$ and $\hat{\tau}$ must be iterated to achieve a maximum. Computation requires solution of simultaneous nonlinear equations via algorithms such as the Newton-Raphson, the Method of Scoring, or the EM algorithm. Jennrich and Schluchter (1985) and Fairclough and Helms (1986) compared performances of the algorithms. Lindstrom and Bates (1988) provided arguments favoring the use of the Newton-Raphson approach. The MIXED procedure in the SAS[®] System software uses a ridge-stabilized Newton-Raphson algorithm and offers, optionally, Method of Scoring steps for the initial iterations. Wolfinger, Tobias, and Sall (1994) reported many further details.

All three algorithms are iterative. Given results from iteration $t \in \{1, 2, \dots\}$, all three algorithms compute a new estimate of the form

$$\hat{\theta}_{t+1} = \hat{\theta}_t + \lambda_t \mathbf{H}_t^{-1} \partial \log L / \partial \theta \Big|_{\theta = \hat{\theta}_t}, \quad (14.11)$$

with $\hat{\theta}_{t+1}$ being the next value of the parameter estimates, $\hat{\theta} = [\hat{\beta}' \ \hat{\tau}']'$. The length of the next step taken towards the MLE is indicated by scalar $\lambda_t \in [0, 1]$. Vector $\mathbf{H}_t^{-1} \partial \log L / \partial \theta$ specifies the direction of step t . In the Newton-Raphson algorithm \mathbf{H}_t is the negative of the observed information matrix,

$$\mathcal{I}(\hat{\theta}) = - \frac{\partial^2}{\partial \theta \partial \theta'} \log L(\theta) \Big|_{\theta = \hat{\theta}}. \quad (14.12)$$

In the Fisher Scoring algorithm \mathbf{H}_t is the negative of the expected information matrix. At convergence, the two algorithms yield a computed value \mathbf{H}_t^{-1} which is an estimate of the asymptotic covariance matrix of $\hat{\theta}_t$. Maximum likelihood theory gives $\mathbf{H}_t^{-1/2} (\hat{\theta}_t - \theta) \overset{\text{asy}}{\approx} \mathcal{N}(\mathbf{0}, \mathbf{I})$. Hence the matrix is useful computing the estimated asymptotic standard errors of $\hat{\theta}_t$, with sufficiently large N .

The EM algorithm has advantages of simplicity, positive definite $\mathbf{H}_t \ \forall t$, and assured increments in the likelihood at each step. However, its iterations provide only the MLEs of $\hat{\beta}$ and τ and do not provide any standard errors. Thus, additional computations are needed to estimate standard errors of $\hat{\beta}$ and $\hat{\tau}$. Of several proposed methods, the supplemental EM algorithm (SEM) of Meng and Rubin (1991) seems best. The SEM derives necessary information about the asymptotic standard errors via evaluation of the rate of convergence of the EM algorithm.

The key parts of conventional maximum likelihood estimation are 1) specification of the likelihood (which arises from assumptions about the model), 2) numerical computation of the MLEs of $\hat{\beta}$ and $\hat{\tau}$, 3) computation of estimates of asymptotically correct approximations of the standard errors of $\hat{\beta}$ and $\hat{\tau}$, and 4) subsequent computation of confidence intervals and test statistics. Under certain regularity conditions the ML estimators have the desirable properties of being consistent, asymptotically Gaussian, and efficient (Harville, 1977; Magnus, 1978). As noted earlier, Kackar and Harville (1984) proved $\hat{\beta}$ is unbiased, while variance estimators are optimistically biased (too small).

14.6 REML ESTIMATION

Restricted maximum likelihood (REML) estimators of variance are less biased because they take into account the loss of degrees of freedom due to the estimation of β . For purposes of estimation (as distinct from hypothesis testing), the method seems preferable to ML methods.

In the simplest $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with N i.i.d. errors, the estimators of β and σ^2 have very desirable properties. With Gaussian data, Lemma 11.6 guarantees the independence of $\tilde{\beta}$ and the residuals \tilde{e} in the $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$. It follows immediately that the ML estimators $\tilde{\beta}$ and $\tilde{\sigma}^2 = \tilde{e}'\tilde{e}/N$ are statistically independent. However, the ML estimator of variance is optimistically small. Hence it is customary to replace the MLE $\tilde{\sigma}^2$ with $\hat{\sigma}^2 = \tilde{\sigma}^2 N/(N - r)$. Here N is the number of ISUs and $r = \text{rank}(\mathbf{X})$. The advantage of $\hat{\sigma}^2$ over $\tilde{\sigma}^2$ lies in the fact that $E(\hat{\sigma}^2) = \sigma^2$, while $E(\tilde{\sigma}^2) = \sigma^2(N - r)/N \leq \sigma^2$.

In seeking estimators with similar properties for mixed models, Patterson and Thompson (1971) recommended transforming the data to functionally separate computation of covariance and expected-value parameter estimators. With

$$\mathbf{y}'_s = [\mathbf{y}'_1 \cdots \mathbf{y}'_N] \tag{14.13}$$

$$\mathbf{X}'_s = [\mathbf{X}'_1 \cdots \mathbf{X}'_N], \tag{14.14}$$

they considered the transformation

$$\hat{\mathbf{e}}_s = [\mathbf{I}_n - \mathbf{X}_s(\mathbf{X}'_s\mathbf{X}_s)^{-1}\mathbf{X}'_s]\mathbf{y}_s. \tag{14.15}$$

As always, $n = \sum_{i=1}^N p_i$ indicates the total number of observations (not the number of ISUs). The appeal of the transformation arises from properties of the $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with i.i.d. Gaussian errors, as summarized in Lemmas 11.5 (independence of $\hat{\mathbf{y}}$ and $\hat{\mathbf{e}}$) and 11.6 (independence of $\tilde{\beta}$ and $\tilde{\mathbf{e}}$). Harville (1974) discussed the concept in the context of estimating variance components. Applying the transformation to the mixed model gives the reduced-profile log likelihood:

$$-2\log L_{\text{REML}}(\hat{\tau}) = (n-q)\log(2\pi) + \sum_{i=1}^N (\log|\Sigma_i| + \hat{e}'_{+i}\Sigma_i^{-1}\hat{e}_{+i}) + \log\left|\sum_{i=1}^N \mathbf{X}'_i\Sigma_i^{-1}\mathbf{X}_i\right|. \tag{14.16}$$

The REML estimates are found by maximizing the last equation. The REML and ML forms differ only by

$$-2[\log L_{\text{REML}}(\hat{\tau}) - \log L_{\text{ML}}(\hat{\tau})] = -q\log(2\pi) + \log\left|\sum_{i=1}^N \mathbf{X}'_i\Sigma_i^{-1}\mathbf{X}_i\right|. \tag{14.17}$$

Lindstrom and Bates (1988) provided helpful discussion about algorithms, as did Demidenko (2004).

14.7 SMALL-SAMPLE PROPERTIES OF ESTIMATORS

Kackar and Harville (1984) proved $\hat{\beta}$ from iterated approximate least squares (sometimes called estimated generalized least squares, among other names) is unbiased. However, covariance parameter estimators, at least in small samples, typically are biased. Littell (2003) provided an excellent overview.

Theorem 14.1 (a) In the $\text{LMM}_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\beta, \mathbf{Z}_i\Sigma_{d_i}(\tau_d)\mathbf{Z}'_i + \Sigma_{e_i}(\tau_e)]$ with Gaussian errors and full-rank $\mathbf{X}_s = [\mathbf{X}'_1 \mathbf{X}'_2 \cdots \mathbf{X}'_N]'$, the ML estimator of β , namely $\hat{\beta}$, is unbiased.

(b) If d_i and e_i have symmetric but not necessarily Gaussian distributions and all other assumptions are met, the result still holds.

(c) If d_i and e_i have symmetric but not necessarily Gaussian distributions, and all other assumptions are met, the REML, method-of-moments (MM) and MINQUE estimators of β are also unbiased.

(d) In general, ML, REML, estimators of τ , and functions thereof, including $\Sigma_{d_i}(\tau_d)$ and $\Sigma_{e_i}(\tau_e)$ are biased.

Proof. A surprisingly simple proof arises from the combination of the symmetry of distribution assumption and a variance estimator expressible as an even function. Demidenko (2004, Section 3.6) provided detailed proofs.

Theorem 14.2 (a) In the *general* $\text{LMM}_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\beta, \mathbf{Z}_i\Sigma_{d_i}(\tau_d)\mathbf{Z}'_i + \Sigma_{e_i}(\tau_e)]$, with or without Gaussian errors, the ML, REML, MM, and MINQUE variance estimators of τ and functions thereof, including $\Sigma_{d_i}(\tau_d)$ and $\Sigma_{e_i}(\tau_e)$, are biased.

(b) In the special case of a balanced random-coefficient model, the REML, MM, and MINQUE variance estimators (1) coincide, (2) are unbiased, and (3) differ from corresponding ML estimators only by a scaling constant, such as $N/(N-1)$.

Proof. Dimedenko's (2004, p. 140) proof of his Theorem 14 and surrounding discussion contain the desired results.

14.8 LARGE-SAMPLE PROPERTIES OF VARIANCE ESTIMATORS

Given the preceding theorem, it is not surprising that only asymptotic results are available for properties of variance parameter estimators $\hat{\tau}$ and $\Sigma_i(\hat{\tau})$. Most such results provide only very low order approximations. As often happens with variance estimation, any bias is nearly always optimistic (estimates too small). In turn, estimated confidence intervals and tests are also optimistic. Much useful work remains to be done.

14.9 CONDITIONAL ESTIMATION OF d_i AND BLUP PREDICTION

In addition to estimation of $\hat{\beta}$ and $\hat{\tau}$, prediction of $\{d_i\}$ (collectively d_s) is often of interest in applications of the general linear mixed model. Henderson (1963) popularized the use of the “best linear unbiased predictor” (BLUP) of d_i , namely $d_i = \Sigma_{di} Z_i' \Sigma_i^{-1} (y_i - X_i \beta)$, with Σ_{di} and Σ_i as in Definition 5.1. For Gaussian data the BLUP is easily proven to be the expected value of d_i conditional on the observed value of y_i . Substituting the MLEs for the unknown parameters yields the empirical BLUP (eBLUP):

$$\hat{d}_i = \hat{\Sigma}_{di} Z_i' \hat{\Sigma}_i^{-1} (y_i - X_i \hat{\beta}). \tag{14.18}$$

Collectively, with Σ_s , Σ_{ds} , and Σ_{es} as defined in equation 5.9,

$$\hat{d}_s = \hat{\Sigma}_{ds} Z_s' \hat{\Sigma}_s^{-1} (y_s - X_s \hat{\beta}). \tag{14.19}$$

Depending on $\hat{\Sigma}_i^{-1}$ and the dimension and magnitude of $(y_i - X_i \hat{\beta})$, the value of \hat{d}_i is subject to shrinkage toward zero. If p_i is small, then conditioning d_i on y_i is not highly informative and the shrinkage toward zero will be substantial. In turn, the eBLUP of the ISU central tendency, namely,

$$\hat{y}_i = X_i \hat{\beta} + Z_i \hat{d}_i \tag{14.20}$$

shrinks toward the estimated population central tendency, $X_i \hat{\beta}$.

Computing the expression for \hat{d}_i above would require inversion of a $p_i \times p_i$ matrix $\hat{\Sigma}_i^{-1}$, which can be problematic for large p_i . Henderson (1963) offered equivalent mixed model equations of the form

$$\begin{bmatrix} X_s' \hat{\Sigma}_{es}^{-1} X_s & X_s' \hat{\Sigma}_{es}^{-1} Z_s \\ Z_s' \hat{\Sigma}_{es}^{-1} X_s & \hat{\Sigma}_{ds}^{-1} + Z_s' \hat{\Sigma}_{es}^{-1} Z_s \end{bmatrix} \begin{bmatrix} \hat{\beta} \\ \hat{d} \end{bmatrix} = \begin{bmatrix} X_s' \hat{\Sigma}_{es}^{-1} y \\ Z_s' \hat{\Sigma}_{es}^{-1} y \end{bmatrix}, \tag{14.21}$$

which involve the more manageable inversion of $\hat{\Sigma}_{es}$ and $\hat{\Sigma}_{ds}$ (as defined in

equation 5.9). The greatest advantage accrues with diagonal $\widehat{\Sigma}_{es}$ and block diagonal $\widehat{\Sigma}_{ds}$.

Harville (1990) proved that

$$\mathcal{V}\left(\begin{bmatrix} \widehat{\beta} \\ \widehat{\mathbf{d}}_s - \mathbf{d}_s \end{bmatrix}\right) = \begin{bmatrix} \mathbf{X}'_s \Sigma_{es}^{-1} \mathbf{X}_s & \mathbf{X}'_s \Sigma_{es}^{-1} \mathbf{Z}_s \\ \mathbf{Z}'_s \Sigma_{es}^{-1} \mathbf{X}_s & \Sigma_{ds}^{-1} + \mathbf{Z}'_s \Sigma_{es}^{-1} \mathbf{Z}_s \end{bmatrix}^{-1}, \quad (14.22)$$

in which both $\widehat{\mathbf{d}}_s$ and \mathbf{d}_s are random vectors. The MLEs must be substituted for the unknown parameters Σ_{es} and Σ_{ds} to use the expression computing estimated standard errors, confidence intervals, and prediction intervals.

EXERCISES

14.1 The LMM $_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i \boldsymbol{\beta}, \mathbf{Z}_i \boldsymbol{\Sigma}_{di}(\boldsymbol{\tau}_d) \mathbf{Z}'_i + \boldsymbol{\Sigma}_{ei}(\boldsymbol{\tau}_e)]$ has $\mathbf{Z}_i = \mathbf{0}$, $p_i = p = 3$, $\boldsymbol{\Sigma}_{ei}(\boldsymbol{\tau}_e) = \begin{bmatrix} \sigma_1^2 & 0 & 0 \\ 0 & \sigma_2^2 & 0 \\ 0 & 0 & \sigma_3^2 \end{bmatrix}$, and $\boldsymbol{\tau} = [\sigma_1^2 \quad \sigma_2^2 \quad \sigma_3^2]'$.

14.1.1 Show that $\boldsymbol{\Sigma}_i(\boldsymbol{\tau})$ has a linear structure.

14.1.2 Assuming $\widehat{\boldsymbol{\tau}}$ is given, provide an appropriate expression for $\widehat{\boldsymbol{\Sigma}}_i$.

14.1.3 Use equation (14.8) to find a slightly simplified expression for $\boldsymbol{\beta}(\boldsymbol{\tau})$.

14.1.4 Use equation (14.10) to provide an expression for $\widehat{\boldsymbol{\tau}} = [\widehat{\sigma}_1^2 \quad \widehat{\sigma}_2^2 \quad \widehat{\sigma}_3^2]'$. Simplify the expression.

14.2 Give an example of a covariance matrix that does *not* have a linear structure.

CHAPTER 15

Tests for Univariate Linear Models

15.1 MOTIVATION

In the present chapter, we always assume Gaussian errors in considering a univariate linear model, a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with $\text{rank}(\mathbf{X}) = r \leq q$, or a GGLM. Throughout the chapter, $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is $a \times 1$ and $\text{rank}(\mathbf{C}) = a \leq q$. With $\boldsymbol{\theta}_0$ an $a \times 1$ vector of known constants, the associated general linear hypothesis (GLH) may be stated

$$H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0, \tag{15.1}$$

with corresponding alternative $H_A : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$. In terms of Boolean variables, with $\mathbb{B}(\cdot) \in \{0, 1\}$, the hypothesis may be written $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$. The first issue addressed is testability. We will restrict attention to well-defined hypotheses and then examine the properties of the resulting tests. We will demonstrate that reasonable hypothesis testing procedures do not exist for nontestable hypotheses.

Another important issue is whether $\boldsymbol{\theta}$ is truly a fixed constant. If the model and its parameters are defined prior to collection of the data, then $\boldsymbol{\theta}$ is an unknown, fixed constant. However, in the course of an analysis one frequently discovers interesting aspects of the data which were not considered prior to data collection and which are directly suggested by the results at hand. In such a case the dimensions and the definition of $\boldsymbol{\theta}$ (implied by the choice of \mathbf{C}) or the model itself may have been influenced by the observed value of \mathbf{y} . Hypotheses suggested by the data (\mathbf{y}) raise issues of multiplicity which may or may not be intractable.

Definition 15.1 (a) Parameters defined without regard to \mathbf{y} are *a priori* parameters.
(b) Other parameters, including ones which are suggested by the data, are *post hoc* parameters.

As indicated by the status of the parameter, the corresponding hypothesis is either an a priori hypothesis or a post hoc hypothesis. Different statistical tests are required for the two different kinds of hypotheses. When appropriate methods are

used, tests of a priori hypotheses are more powerful than tests of post hoc hypotheses.

15.2 TESTABILITY OF UNIVARIATE HYPOTHESES

In the theoretical formulation of any hypothesis test procedure, the hypothesis is always manipulated in terms of its linearly independent (LIN) components. Hence one requires the rows of \mathbf{C} to be LIN, with $\text{rank}(\mathbf{C}) = a \leq q$. In the following developments, we first address the issue of testability for full-rank \mathbf{C} matrices. Later we consider less-than-full-rank \mathbf{C} . With full-rank \mathbf{C} , the formal definition of a testable hypothesis implies questions of testability that arise only in LTFR models.

Definition 15.2 Under the assumptions of $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with $\text{rank}(\mathbf{X}) = r \leq q$ and $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ $a \times 1$, the hypothesis $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$ is *testable* if and only if $\boldsymbol{\theta}$ is estimable and $\text{rank}(\mathbf{C}) = a \leq q$.

In FR models, $\boldsymbol{\beta}$ and $\boldsymbol{\theta}$ are always estimable and corresponding tests with full-rank \mathbf{C} are always testable. In LTFR models $H_0 = \mathbb{B}(\boldsymbol{\beta} = \mathbf{0})$ versus $H_A = \mathbb{B}(\boldsymbol{\beta} \neq \mathbf{0})$ is never testable, while $\boldsymbol{\theta}$ may or may not be estimable and hence may or may not be testable.

The great majority of all results presented about properties of hypothesis tests depend on the assumption of Gaussian errors. In contrast, in estimation theory, many first- and second-moment properties are distribution free in the sense that the particular likelihood need not be specified (finite second moments and appropriate independence and homogeneity suffice). The distribution-free feature carries over to testability, even though the forms of the distributions of test statistics studied here depend strongly on the Gaussian assumption.

The following argument helps explain the relationship between estimable parameters and testable hypotheses. We will find that test statistics are functions of $\boldsymbol{\theta} = \mathbf{C}\tilde{\boldsymbol{\beta}}$, with $\tilde{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}$ in the FR case. The quantity $\tilde{\boldsymbol{\theta}}$ needs to be invariant to the choice of $\tilde{\boldsymbol{\beta}}$ in order for the test to be invariant. We do not want the test result to depend on our choice of generalized inverse! We have previously proven $\tilde{\boldsymbol{\theta}}$ is invariant if and only if $\boldsymbol{\theta}$ is estimable. A series of theorems were presented in Chapter 10 which give estimability criteria. In particular, $\mathbf{C}(\mathbf{X}'\mathbf{X})^-(\mathbf{X}'\mathbf{X}) = \mathbf{C}$ if and only if $\mathbf{C}\boldsymbol{\beta}$ is estimable.

Lemma 15.1 For $a \times q$ \mathbf{C} , the matrix $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^-\mathbf{C}'$ is $a \times a$.

- (a) If \mathbf{X} is full rank, then $\text{rank}(\mathbf{M}) = a$ provides a necessary and sufficient condition for testability.
- (b) If \mathbf{X} is less than full rank, then the condition $\text{rank}(\mathbf{M}) = a$ provides a necessary but not sufficient condition for testability.

(e) If X is less than full rank, then $\text{rank}(M) = a$ combined with the requirement $C = C(X'X)^-(X'X)$, or any other condition which guarantees estimability, provide a necessary and sufficient set of conditions for testability.

Proof. Left as an exercise.

Lemma 15.2 Estimable $\theta = C\beta$ ($a \times 1$) with $\text{rank}(C) = a$ implies M is symmetric, and unique and $\text{rank}(M) = a$. Also $C = AX$ with $\text{rank}(A) = a$.

Proof. Estimability allows writing $C = AX$ (Theorem 11.3). In turn,

$$\begin{aligned} M &= C(X'X)^-C' \\ &= AX(X'X)^-X'A' \\ &= AX(X'X)^+X'A' \\ &= AHA'. \end{aligned} \tag{15.2}$$

The matrix $H = X(X'X)^+X'$ is symmetric, idempotent, unique and of rank $r = \text{rank}(X)$ (Theorem 1.15). Hence M is symmetric whether or not $(X'X)^-$ is symmetric, and M is unique.

Lemma 1.29 gives $X = L_1 \text{Dg}(s_1) R_1'$, with $R_1' R_1 = I_r$, $\text{Dg}(s_1)$ $r \times r$ of rank r , L_1 $N \times r$ of rank r , $L_1' L_1 = I_r$, and $H = L_1 L_1'$. Hence $M = (AL_1)(AL_1)'$. Having $C = AX$ implies $CC' = AXX'A' = AL_1 \text{Dg}(s_1)^2 L_1' A' = FF'$ for $F = [AL_1 \text{Dg}(s_1)]$ and $a = \text{rank}(C) = \text{rank}(CC') = \text{rank}(F) = \text{rank}(AL_1) = \text{rank}(M)$. With $C = AX$ of rank $a \leq r$, Lemma 1.6 gives $\text{rank}(AX) \leq \min[\text{rank}(A), \text{rank}(X)]$, which implies $a \leq \min[\text{rank}(A), r]$ and $a \leq \text{rank}(A)$. Also A is $a \times N$, with $a \leq N$, which implies $\text{rank}(A) \leq a$. Hence $\text{rank}(A) = a$. \square

Example 15.1 The theory of linearly equivalent models provides many insights into the structure underlying a testable hypothesis. For any testable $a \times 1$ $\theta = C\beta$, the last lemma and the notation developed in its proof allow defining full-rank $N \times r$ matrix $X_1 = XR_1 = L_1 \text{Dg}(s_1)$ and $r \times 1$ matrix $\beta_1 = R_1' \beta$. The corresponding model equation may be written

$$\begin{aligned} y &= X \beta + e \\ &= L_1 \text{Dg}(s_1) R_1' \beta + e \\ &= X_1 \beta_1 + e. \end{aligned} \tag{15.3}$$

Choosing the $a \times r$ matrix $C_1 = CR_1$ implies $C_1 R_1' = C$ and $\theta_1 = C_1 \beta_1 = C R_1' \beta = C\beta$. Orthonormal P and Q ensure $(PAQ')^+ = QA^+P'$. Also

$$\begin{aligned} M_1 &= C_1(X_1'X_1)^-C_1' = CR_1(X_1'X_1)^-R_1'C' \\ &= C(R_1X_1'X_1R_1')^+C' \\ &= M. \end{aligned} \tag{15.4}$$

Hence with a less-than-full-rank model, a testable hypothesis corresponds to finding the linearly equivalent full-rank model with predictors $X_1 = XR_1$, the

principal component scores. In turn, C_1 expresses the original contrasts in terms of the component score parameters, namely β_1 .

Theorem 15.1 For $GLM_{N,q}(y_i; X_i\beta, \sigma^2)$ with $\text{rank}(X) = r \leq q$, $\theta = C\beta$ $a \times 1$, and $\text{rank}(C) = a \leq q$, if θ is *not* estimable then $H_0 = \mathbb{B}(\theta = \theta_0)$ versus $H_A = \mathbb{B}(\theta \neq \theta_0)$ is not testable. Consequently no reasonable test exists for the hypothesis; i.e., any test procedure will have size $\alpha = 0$ and power $1 - \beta = 0$. Furthermore the test will not be invariant to changes in irrelevant quantities.

Proof. (Searle, 1971, p. 193–194) Intuitively, tests involve a comparison of the model with and without the linear restrictions $C\beta = \theta_0$. In the case of nontestable hypotheses the model seems to fit equally well with or without the restrictions. If (a) θ is not estimable, (b) $\tilde{\beta}$ satisfies the restricted normal equations $X'X\tilde{\beta} + C'\lambda = X'y$, and (c) $C\tilde{\beta} = \theta_0$ (in which λ is the vector of “Lagrange multipliers”), then $\tilde{\beta}$ also satisfies the unrestricted normal equations $X'X\tilde{\beta} = X'y$. For a nontestable hypothesis, *SSE* is the same with or without restrictions! In terms of goodness of fit, the unrestricted model and the restricted model are indistinguishable; we have no reason to prefer one to the other. Thus, any test has size $\alpha \equiv \text{Pr}\{\text{rejecting } H_0 | H_0 = \text{TRUE}\} = 0$ and power $1 - \beta \equiv \text{Pr}\{\text{reject } H_0 | H_A = \text{TRUE}\} = 0$. \square

Corollary 15.1 If $\theta_s = [\theta'_1 \ \theta'_2]'$ and θ_1 is estimable but θ_2 is not, then tests for θ_s are indistinguishable from tests for θ_1 .

Proof. For $\theta_0 = [\theta'_{0,1} \ \theta'_{0,2}]'$ the model subject to the restrictions $C_1\beta = \theta_{0,1}$ is linearly equivalent to some unrestricted model. It then follows from the proof for Theorem 11.18 (with restrictions not necessarily of full rank) that, in terms of goodness of fit via *SSE*, the model constrained by both $\theta_1 \equiv C_1\beta - \theta_{0,1} = 0$ and $\theta_2 \equiv C_2\beta - \theta_{0,2} = 0$ cannot be distinguished from the model constrained only by $\theta_1 \equiv C_1\beta - \theta_{0,1} = 0$. Therefore either of the two constrained models may be compared with the original, unrestricted model; the results in terms of difference in *SSE* would be the same. \square

Test statistics are functions of $\hat{\delta} = (C\tilde{\beta} - \theta_0)'[C(X'X)^{-1}C']^{-1}(C\tilde{\beta} - \theta_0)$ (with $\tilde{\beta} = \hat{\beta}$ and $X'X$ nonsingular in the FR case) and the test statistics depend on C and θ_0 only through $\hat{\delta}$. The ability to compute $\hat{\delta}$ depends on the existence of $M^{-1} = [C(X'X)^{-1}C']^{-1}$ and can be computed if and only if $M = [C(X'X)^{-1}C'] \neq 0$ has full rank. Having testable θ ensures full-rank M .

In many cases it is possible to compute $\hat{\delta}$ (and therefore the test statistic) even though θ is *not* estimable and the hypothesis is *not* testable. The result is formalized in the following theorem. In turn, the next theorem answers the question “What hypothesis is the test statistic addressing?” in such cases.

Theorem 15.2 If $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ has $\text{rank}(\mathbf{X}) = r \leq q$, $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ ($a \times 1$), and $\text{rank}(\mathbf{C}) = a \leq q$, then

- (a) $\boldsymbol{\theta}$ is estimable implies $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}'$ has full rank, although
- (b) $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}'$ having full rank does *not* imply $\boldsymbol{\theta}$ is estimable.

Proof of (a). (Compare with Searle, 1971, p. 189–190) It is assumed $a \times q$ \mathbf{C} has (full) $\text{rank}(\mathbf{C}) = a \leq q$. Estimable $\mathbf{C}\boldsymbol{\beta}$ insures there exists unique matrix \mathbf{A} such that $\mathbf{C} = \mathbf{A}(\mathbf{X}'\mathbf{X})$ which implies $\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}' = \mathbf{A}(\mathbf{X}'\mathbf{X})\mathbf{A}'$. Therefore $\text{rank}[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}'] = \text{rank}[(\mathbf{A}\mathbf{X}')(\mathbf{X}\mathbf{A}')] \equiv \text{rank}(\mathbf{A}\mathbf{X}') \leq a$.

It suffices to prove $\text{rank}(\mathbf{A}\mathbf{X}') = a$. Matrix theory gives $a \equiv \text{rank}(\mathbf{C}) \equiv \text{rank}[\mathbf{A}(\mathbf{X}'\mathbf{X})] \leq \min\{\text{rank}(\mathbf{A}\mathbf{X}'), \text{rank}(\mathbf{X})\}$. Therefore $\text{rank}(\mathbf{A}\mathbf{X}') \geq a$ and $\text{rank}(\mathbf{X}) \equiv r \geq a$. Hence $\text{rank}(\mathbf{A}\mathbf{X}') = a$ since $\text{rank}(\mathbf{A}\mathbf{X}')$ cannot exceed a . We can similarly prove $\text{rank}(\mathbf{A}) = a$ since $a \equiv \text{rank}(\mathbf{C}) \equiv \text{rank}[\mathbf{A}(\mathbf{X}'\mathbf{X})] \leq \min\{\text{rank}(\mathbf{A}), \text{rank}(\mathbf{X}'\mathbf{X})\}$.

Proof of (b). By counterexample, with $p = 1, q = N = 3$:

$$\mathbf{X}\boldsymbol{\beta} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{bmatrix} \tag{15.5}$$

$$(\mathbf{X}'\mathbf{X})^{-} = \left(\begin{bmatrix} 3 & 2 & 1 \\ 2 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \right)^{-} = \frac{1}{2} \begin{bmatrix} 2 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix} \tag{15.6}$$

$$\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} \mu \\ \alpha_1 \\ \alpha_2 \end{bmatrix} = \begin{bmatrix} \mu \\ \alpha_1 \end{bmatrix} \tag{15.7}$$

$$\mathbf{M}^{-1} = \begin{bmatrix} 3 & 2 \\ 2 & 2 \end{bmatrix} \tag{15.8}$$

$$\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}(\mathbf{X}'\mathbf{X}) = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \end{bmatrix} \neq \mathbf{C}. \tag{15.9}$$

Hence the result is true by counterexample. □

Theorem 15.3 For $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with $\text{rank}(\mathbf{X}) = r \leq q$, $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ $a \times 1$, and $\text{rank}(\mathbf{C}) = a \leq q$, even though $\mathbf{M}^{-1} = [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}']^{-1}$ exists, $\boldsymbol{\theta}$ may not be estimable. If so, test statistics based on $\hat{\delta} = (\mathbf{C}\tilde{\boldsymbol{\beta}} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1} (\mathbf{C}\tilde{\boldsymbol{\beta}} - \boldsymbol{\theta}_0)$ can be computed and the testable hypothesis actually tested is $H_0 = \mathbb{B}(\boldsymbol{\theta}_G = \mathbf{0})$ versus $H_A = \mathbb{B}(\boldsymbol{\theta}_G \neq \mathbf{0})$ with $\boldsymbol{\theta}_G = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}(\mathbf{X}'\mathbf{X})\boldsymbol{\beta} - \boldsymbol{\theta}_0$.

Proof. Searle (1971, p. 195) provided a proof.

Here θ_G is not necessarily invariant to the choice of $(\mathbf{X}'\mathbf{X})^-$. For each choice of generalized inverse a potentially different testable hypothesis is tested! The special case of $r = q$ has all θ estimable. In the special case, full rank of $\mathbf{M} = [\mathbf{C}(\mathbf{X}'\mathbf{X})^- \mathbf{C}']$ guarantees a testable hypothesis, which corresponds to a unique parameter and a well-defined test.

Example 15.2 The following example illustrates a poorly defined hypothesis. If $p = 1$, $q = N = 3$, $\theta_0 = \mathbf{0}$,

$$\mathbf{X} = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \quad (15.10)$$

$$(\mathbf{X}'\mathbf{X})^- = \left(\begin{bmatrix} 3 & 2 & 1 \\ 2 & 2 & 0 \\ 1 & 0 & 1 \end{bmatrix} \right)^- = \frac{1}{2} \begin{bmatrix} 2 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix} \quad (15.11)$$

$$\mathbf{C} = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad (15.12)$$

$$\mathbf{C}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'\mathbf{X} = \mathbf{C} \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{bmatrix}, \quad (15.13)$$

$\boldsymbol{\beta} = [\mu \ \alpha_1 \ \alpha_2]'$, $\boldsymbol{\theta} = [\mu \ \alpha_1]'$, and $\boldsymbol{\theta}_G = [\mu + \alpha_2 \ \alpha_1 - \alpha_2]'$. More generally, all possible generalized inverses are given by

$$(\mathbf{X}'\mathbf{X})^- = \frac{1}{2} \begin{bmatrix} 2 & -2 & 0 \\ -2 & 3 & 0 \\ 0 & 0 & 0 \end{bmatrix} + \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{bmatrix} \mathbf{T} + \mathbf{S} \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ -1 & 1 & 1 \end{bmatrix}, \quad (15.14)$$

in which \mathbf{T} and \mathbf{S} are arbitrary. Choosing

$$(\mathbf{X}'\mathbf{X})^- = \frac{1}{2} \begin{bmatrix} 2 & -2 & -2 \\ -2 & 3 & 2 \\ 0 & 0 & 2 \end{bmatrix} \quad (15.15)$$

would give $\boldsymbol{\theta}_G = [0 \ \mu + \alpha_1]'$.

We now consider more general forms of \mathbf{C} . Previously we have assumed $a \times q$ \mathbf{C} with $1 \leq \text{rank}(\mathbf{C}) = a \leq q$. Can the requirements be relaxed? We are occasionally interested in less-than-full-rank \mathbf{C} . In any case we will need $\mathbf{C} = \mathbf{A}\mathbf{X}$ (for some \mathbf{A}) to have any hope of formulating a testable hypothesis. Otherwise $\boldsymbol{\theta}$ will not be estimable. Estimability of $\boldsymbol{\theta}$ does not require \mathbf{C} to have full rank. In particular, if $\mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}$ is estimable, then so is

$$\begin{bmatrix} \mathbf{C} \\ \mathbf{C} \end{bmatrix} \boldsymbol{\beta} = \mathbf{C}_2 \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\theta} \\ \boldsymbol{\theta} \end{bmatrix} = \boldsymbol{\theta}_2. \tag{15.16}$$

Here $1 \leq \text{rank}(\mathbf{C}_2) < 2a$. Also, $H_0 = \mathbb{B}(\boldsymbol{\theta}_2 = \mathbf{0}) = \mathbb{B}[(\boldsymbol{\theta} = \mathbf{0}) \cap (\boldsymbol{\theta} = \mathbf{0})]$ is overstated and is identical to $H_{0,1} = \mathbb{B}(\boldsymbol{\theta} = \mathbf{0})$. We should test subhypothesis $H_{0,1}$ rather than H_0 itself.

Theorem 15.4 A test with a LTFR \mathbf{C} matrix is indistinguishable from a test with \mathbf{C} replaced by a full (row) rank matrix with rows that span all rows of \mathbf{C} . More specifically, model $\text{GLM}_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ has $\text{rank}(\mathbf{X}) = r \leq q$, estimable $(a \times 1)$

$$\boldsymbol{\theta} = \mathbf{C} \boldsymbol{\beta} = \begin{bmatrix} \mathbf{C}_1 \\ \mathbf{C}_2 \end{bmatrix} \boldsymbol{\beta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix}, \tag{15.17}$$

and $\boldsymbol{\theta}_0 = [\boldsymbol{\theta}'_{0,1} \ \boldsymbol{\theta}'_{0,2}]'$. Also \mathbf{C}_j , $\boldsymbol{\theta}_j$, and $\boldsymbol{\theta}_{0,j}$ have a_j rows for $j \in \{1, 2\}$, $\text{rank}(\mathbf{C}) = \text{rank}(\mathbf{C}_1) = a_1 < \min\{a, q\}$.

(a) Tests about $\boldsymbol{\theta}$ are indistinguishable from tests about $\boldsymbol{\theta}_1$; i.e., $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ iff $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{0,1}$, and $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0) = \mathbb{B}(\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{0,1})$.

(b) Tests about $\boldsymbol{\theta}_1 = \mathbf{C}_1 \boldsymbol{\beta}$ are testable; i.e., \mathbf{C}_1 has full row rank and $\boldsymbol{\theta}_1$ is estimable.

Proof. The rows of \mathbf{C}_2 are linear combinations of the rows of \mathbf{C}_1 ; i.e., $\mathbf{C}_2 = \mathbf{A} \mathbf{C}_1$ for some fixed matrix \mathbf{A} . Hence $\boldsymbol{\theta}_2 = \mathbf{A} \mathbf{C}_1 \boldsymbol{\beta} = \mathbf{A} \boldsymbol{\theta}_1$. Also, $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ iff $\boldsymbol{\theta}_1 = \boldsymbol{\theta}_{0,1}$ and

$$\mathbf{C} \boldsymbol{\beta} = \begin{bmatrix} \mathbf{I} \\ \mathbf{A} \end{bmatrix} \mathbf{C}_1 \boldsymbol{\beta} = \boldsymbol{\theta} = \begin{bmatrix} \boldsymbol{\theta}_1 \\ \boldsymbol{\theta}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{I} \\ \mathbf{A} \end{bmatrix} \boldsymbol{\theta}_1. \tag{15.18}$$

Estimable $\boldsymbol{\theta}$ implies $\boldsymbol{\theta}_1$ is estimable; combining estimability with \mathbf{C}_1 having full row rank implies $\widehat{\boldsymbol{\delta}}_1$ exists. While $\widehat{\boldsymbol{\delta}} = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}']^{-1} (\boldsymbol{\theta} - \boldsymbol{\theta}_0)$ does not exist and cannot be computed, $\widehat{\boldsymbol{\delta}}_1 = (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_{0,1})' [\mathbf{C}_1(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_1]^{-1} (\boldsymbol{\theta}_1 - \boldsymbol{\theta}_{0,1})$ does exist and can be computed. \square

The following theorem justifies using a simpler model which is linearly equivalent. The approach allows avoiding some of the pitfalls with LTFR models.

Theorem 15.5 Any primary or secondary expected-value parameter testable in $\text{GLM}_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ is also testable in a linearly equivalent model.

Proof. Follows almost directly from the parallel result about estimability.

15.3 TESTS OF A PRIORI HYPOTHESES

Theorem 15.6 The likelihood ratio test in the univariate full rank GLM may be expressed exactly in terms of an F . For a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ with

Gaussian errors, interest centers on a priori and estimable $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ ($a \times 1$) with $\text{rank}(\mathbf{C}) = a \leq q$, $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$, and \mathbf{C} and $\boldsymbol{\theta}_0$ known constants.

(a) For testing $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$, with $\Pr\{F(\nu_1, \nu_2, \omega) > f_{\text{crit}}\} = \alpha$, the *likelihood ratio test* (LRT) of size α is $\phi(\mathbf{y}) = \mathbb{B}[F(\mathbf{y}) > f_{\text{crit}}]$, in which

$$\widehat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (15.19)$$

$$\widehat{\boldsymbol{\theta}} = \mathbf{C}\widehat{\boldsymbol{\beta}} \quad (15.20)$$

$$\widehat{\sigma}^2 = \mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}/(N - q) \quad (15.21)$$

$$F(\mathbf{y}) = [(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/a] / \widehat{\sigma}^2. \quad (15.22)$$

(b) Under H_0 $F(\mathbf{y}) \sim F(a, N - q, 0)$.

(c) In general, $F(\mathbf{y}) \sim F(a, N - q, \omega)$, with

$$\omega = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) / \sigma^2. \quad (15.23)$$

(d) Testing $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$ is equivalent to testing $H_0 = \mathbb{B}(\omega = 0)$ versus $H_A = \mathbb{B}(\omega \neq 0)$.

(e) Test statistic $F(\mathbf{y})$ can be written as a ratio of chi-square statistics (scaled sums of squares) divided by their degrees of freedom,

$$F(\mathbf{y}) = \frac{SSH/a}{SSE/(N - q)}, \quad (15.24)$$

with SSH = sum of squares for the hypothesis and SSE = sum of squares for error.

Proof. Results in Chapter 11 on univariate linear model estimation provide the basis for the proof, which is cast in terms of the more general concept of supremum, rather than maximum. Doing so simplifies generalizing the proof to LTFR models because FR models lead to a unique estimator of $\boldsymbol{\beta}$, while LTFR models do not.

Part 1. The model has primary parameters $\boldsymbol{\tau} = [\boldsymbol{\beta}' \sigma^2]'$. Unrestricted estimation finds values in the set $\boldsymbol{\tau}_A = \{\boldsymbol{\tau} : \boldsymbol{\beta} \in \mathbb{R}^q, \sigma^2 \geq 0\}$. The null hypothesis restricts attention to the set $\boldsymbol{\tau}_0 = \{\boldsymbol{\tau} : \boldsymbol{\theta} = \boldsymbol{\theta}_0, \boldsymbol{\beta} \in \mathbb{R}^q, \sigma^2 \geq 0\} \subset \boldsymbol{\tau}_A$.

Unrestricted maximum likelihood estimators are derived as follows. The maximum likelihood estimate is any particular value of $\boldsymbol{\tau}$, say $\widehat{\boldsymbol{\tau}}$, for which the likelihood

$$L(\boldsymbol{\tau}; \mathbf{y}_*) = (2\pi\sigma^2)^{-N/2} \exp[-\sigma^{-2}(\mathbf{y}_* - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y}_* - \mathbf{X}\boldsymbol{\beta})/2], \quad (15.25)$$

or equivalently

$$\log L(\boldsymbol{\tau}; \mathbf{y}_*) = -N \log(2\pi) - N \log(\sigma^2) - \sigma^{-2}(\mathbf{y}_* - \mathbf{X}\boldsymbol{\beta})'(\mathbf{y}_* - \mathbf{X}\boldsymbol{\beta})/2, \quad (15.26)$$

achieves its supremum. The LRT statistic is

$$\gamma(\mathbf{y}) = \frac{\sup_{\tau \in \tau_0} L(\boldsymbol{\beta}, \sigma^2; \mathbf{y}_*)}{\sup_{\tau \in \tau_A} L(\boldsymbol{\beta}, \sigma^2; \mathbf{y}_*)} = \frac{L(\hat{\tau}_0)}{L(\hat{\tau}_A)}. \tag{15.27}$$

Here $L(\hat{\tau}_0)$ is the restricted supremum of L over all τ which satisfy H_0 , while $L(\hat{\tau}_A)$ is the unrestricted supremum of L over all τ which satisfy the model.

For the full-rank case, the unrestricted supremum is obtained at $\hat{\tau}_A$ specified by

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}_* \tag{15.28}$$

$$\hat{\sigma}^2 = (\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})/N. \tag{15.29}$$

Hence

$$\begin{aligned} L(\hat{\tau}_A) &= (2\pi\hat{\sigma}^2)^{-N/2} \exp\left[-(\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})/(2\hat{\sigma}^2)\right] \\ &= (2\pi\hat{\sigma}^2)^{-N/2} \exp\left(-\frac{N}{2} \frac{\hat{\sigma}^2}{\hat{\sigma}^2}\right) \\ &= (2\pi\hat{\sigma}^2 \mathbf{e})^{-N/2}. \end{aligned} \tag{15.30}$$

Part 2. The technique of Lagrangian multipliers allows finding the supremum of $\log L(\tau; \mathbf{y}_*)$ subject to the restrictions of $H_0 : \boldsymbol{\theta} = \mathbf{0}$. In the following, $\boldsymbol{\lambda}$ is an $a \times 1$ vector of Lagrangian multipliers. The restricted optimization with respect to τ is achieved by undertaking unrestricted optimization with respect to $[\tau' \boldsymbol{\lambda}']'$ for the objective function

$$h(\tau, \boldsymbol{\lambda}; \mathbf{y}_*) = \log L(\tau; \mathbf{y}_*) - (\mathbf{C}\boldsymbol{\beta} - \boldsymbol{\theta}_0)' \boldsymbol{\lambda}. \tag{15.31}$$

Here $h() = \log L()$ for all values of $\boldsymbol{\beta}$ satisfying the restriction. For taking derivatives with respect to $\boldsymbol{\beta}$ it is useful to write

$$\partial h / \partial \boldsymbol{\lambda} = \mathbf{C}\boldsymbol{\beta} - \boldsymbol{\theta}_0 \tag{15.32}$$

$$\partial h / \partial \boldsymbol{\beta} = \frac{2}{2\sigma^2} \mathbf{X}'\mathbf{y}_* - \frac{2}{2\sigma^2} (\mathbf{X}'\mathbf{X})\boldsymbol{\beta} - \mathbf{C}'\boldsymbol{\lambda} \tag{15.33}$$

$$\partial h / \partial \sigma^2 = \frac{-N}{2\sigma^2} + \frac{1}{2\sigma^4} (\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}}). \tag{15.34}$$

Setting each derivative to zero and simplifying gives

$$\mathbf{C}\bar{\boldsymbol{\beta}} = \boldsymbol{\theta}_0 \tag{15.35}$$

$$\mathbf{X}'\mathbf{X}\bar{\boldsymbol{\beta}} + \mathbf{C}'\boldsymbol{\lambda}\sigma^2 = \mathbf{X}'\mathbf{y}_* \tag{15.36}$$

$$\bar{\sigma}^2 = (\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})/N. \tag{15.37}$$

The second equation implies

$$\bar{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'\boldsymbol{\lambda}\sigma^2. \tag{15.38}$$

Combining the last result with the first equation gives

$$\theta_0 = \mathbf{C}\bar{\boldsymbol{\beta}} = \mathbf{C}\hat{\boldsymbol{\beta}} - \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'\lambda\sigma^2. \quad (15.39)$$

Hence

$$\lambda = [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\mathbf{C}\hat{\boldsymbol{\beta}} - \theta_0)/\sigma^2. \quad (15.40)$$

The last two equations together give

$$\bar{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\mathbf{C}\hat{\boldsymbol{\beta}} - \theta_0). \quad (15.41)$$

Hence $\bar{\boldsymbol{\beta}}$ satisfies the restriction $\mathbf{C}\bar{\boldsymbol{\beta}} = \theta_0$ and is in fact the restricted MLE of $\boldsymbol{\beta}$.

Part 3. Demonstrating $\bar{\boldsymbol{\beta}}$ is the supremum of the function, as desired, may be approached in various ways. A direct approach would be to begin by applying Theorem 9.15 in Schott (2005), which allows verifying a local maximum was achieved. Complete verification also requires excluding boundary values of the parameter space, the set over which maximization was performed, as possible solutions. Alternately, monotonicity properties of the likelihood could be used. Finally, results for the restricted model could be cast in terms of a linearly equivalent unrestricted model.

Part 4. The last form allows writing the restricted maximum of the likelihood as

$$\begin{aligned} L(\hat{\boldsymbol{\tau}}_0) &= (2\pi\bar{\sigma}^2)^{-N/2} \exp[-(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})/(2\bar{\sigma}^2)] \\ &= (2\pi\bar{\sigma}^2)^{-N/2} \exp\left(\frac{-N}{2} \frac{\bar{\sigma}^2}{\bar{\sigma}^2}\right) \\ &= (2\pi\bar{\sigma}^2\mathbf{e})^{-N/2}. \end{aligned} \quad (15.42)$$

In turn, the LRT statistic is

$$\begin{aligned} \gamma(\mathbf{y}) &= \frac{L(\hat{\boldsymbol{\tau}}_0)}{L(\hat{\boldsymbol{\tau}}_A)} = \frac{[(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\bar{\boldsymbol{\beta}})]^{-N/2}}{[(\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})'(\mathbf{y}_* - \mathbf{X}\hat{\boldsymbol{\beta}})]^{-N/2}} \\ &= (\bar{\sigma}^2/\hat{\sigma}^2)^{-N/2}. \end{aligned} \quad (15.43)$$

Part 5. The fact that $0 \leq L(\hat{\boldsymbol{\tau}}_0) \leq L(\hat{\boldsymbol{\tau}}_A)$ allows concluding $\gamma(\mathbf{y}) \in [0, 1]$. If $\|\mathbf{C}\boldsymbol{\beta} - \theta_0\| = (\mathbf{C}\boldsymbol{\beta} - \theta_0)'(\mathbf{C}\boldsymbol{\beta} - \theta_0) \approx 0$, then $L(\hat{\boldsymbol{\tau}}_0) \approx L(\hat{\boldsymbol{\tau}}_A)$ and $\gamma \approx 1$, while if $\|\mathbf{C}\boldsymbol{\beta} - \theta_0\| \gg 0$, then $L(\hat{\boldsymbol{\tau}}_0) \ll L(\hat{\boldsymbol{\tau}}_A)$ and $\gamma(\mathbf{y}) \ll 1$. Hence reject H_0 for small values of $\gamma(\mathbf{y})$ with the decision function

$$\phi_{\text{LRT}}(\mathbf{y}) = \mathbb{B}[\gamma(\mathbf{y}) < \gamma_\alpha]. \quad (15.44)$$

We indicate the density of $\gamma(\mathbf{y})$ by $g(\gamma; \boldsymbol{\tau})$ [which exists due to $\gamma(\mathbf{y})$ being a well-behaved and smooth function of a.c. random variables]. For tests of size α , the appropriate critical value is γ_α as specified by

$$\int_0^{\gamma_\alpha} g(\gamma; \boldsymbol{\tau}|\boldsymbol{\theta} = \boldsymbol{\theta}_0) d\gamma = \alpha. \quad (15.45)$$

The critical region is $[0, \gamma_\alpha]$.

Part 6. The complexity of the density, $g(\gamma)$, led statisticians to work with a different statistic. Using the fact

$$\bar{\beta} = \hat{\beta} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\mathbf{C}\hat{\beta} - \theta_0) \quad (15.46)$$

allows writing

$$\mathbf{y}_* - \mathbf{X}\bar{\beta} = (\mathbf{y}_* - \mathbf{X}\hat{\beta}) + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\hat{\theta} - \theta_0). \quad (15.47)$$

Hence

$$\begin{aligned} (\mathbf{y}_* - \mathbf{X}\bar{\beta})'(\mathbf{y}_* - \mathbf{X}\bar{\beta}) &= (\mathbf{y}_* - \mathbf{X}\hat{\beta})'(\mathbf{y}_* - \mathbf{X}\hat{\beta}) + \\ &\quad (\hat{\theta} - \theta_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\hat{\theta} - \theta_0) \\ &= SSE + SSH, \end{aligned} \quad (15.48)$$

in which $\hat{\theta} = \mathbf{C}\hat{\beta}$,

$$\begin{aligned} SSE &= (\mathbf{y}_* - \mathbf{X}\hat{\beta})'(\mathbf{y}_* - \mathbf{X}\hat{\beta}) \\ &= \mathbf{y}_*'[I - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}_*, \end{aligned} \quad (15.49)$$

and

$$SSH = (\hat{\theta} - \theta_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\hat{\theta} - \theta_0). \quad (15.50)$$

Thus

$$\begin{aligned} \gamma(\mathbf{y}) &= \left[\frac{(\mathbf{y}_* - \mathbf{X}\bar{\beta})'(\mathbf{y}_* - \mathbf{X}\bar{\beta})}{(\mathbf{y}_* - \mathbf{X}\hat{\beta})'(\mathbf{y}_* - \mathbf{X}\hat{\beta})} \right]^{-N/2} \\ &= \left(\frac{SSE + SSH}{SSE} \right)^{-N/2} \\ &= \left(1 + \frac{SSH}{SSE} \right)^{-N/2}. \end{aligned} \quad (15.51)$$

Part 7. Statistical independence of $\hat{\beta}$ and $(\mathbf{y} - \mathbf{X}\hat{\beta})$ guarantees independence of quadratic forms SSE and SSH . Furthermore $SSE/\sigma^2 \sim \chi^2(N - q)$ and $SSH/\sigma^2 \sim \chi^2(a, \omega)$. The properties lead to the transformation

$$\begin{aligned} F &= \frac{SSH/\text{rank}(\mathbf{C})}{SSE/(N - q)} \\ &= \frac{N - q}{a} (\gamma^{-2/N} - 1) \\ &\sim F(a, N - q, \omega), \end{aligned} \quad (15.52)$$

with $\omega = (\theta - \theta_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\theta - \theta_0)/\sigma^2$. Testing $H_0 : \theta = \theta_0$ versus $H_A : \theta \neq \theta_0$ is equivalent to testing $H_0 : \omega = 0$ versus $H_A : \omega \neq 0$ because $\omega = 0$ iff $\theta = \theta_0$.

Here F is a monotone decreasing function of γ and $\gamma \in (0, \gamma_m)$ iff $F \in [F_m, \infty)$ with $F_m = (\gamma_m^{-2/N} - 1)(N - q)/a$. Therefore a test based on F can be made equivalent to the test based on $\gamma : \phi_{\text{LRT}}(\mathbf{y}) = \mathbb{B}[\gamma(\mathbf{y}) < \gamma_\alpha] = \mathbb{B}[F(\mathbf{y}) > f_{\text{crit}}]$. Hypothesis $\boldsymbol{\theta} = \boldsymbol{\theta}_0$ or equivalently $\omega = 0$ is rejected for improbably small values of the likelihood ratio, or equivalently, for improbably large values of the F statistic. \square

Corollary 15.6.1 For $\text{GLM}_{N,q}\text{LTFR}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors and $\text{rank}(\mathbf{X}) = r < q$, $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ ($a \times 1$) with known constant \mathbf{C} of rank $a \leq r$. If $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$ is a testable hypothesis, then the likelihood ratio test may be implemented exactly by rejecting H_0 for improbably large values of the statistic

$$F = \left\{ (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/a \right\} / \hat{\sigma}^2. \quad (15.53)$$

Here $\hat{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}}$, and $\tilde{\boldsymbol{\beta}}$ is any least squares estimator of $\boldsymbol{\beta}$ based on $(\mathbf{X}'\mathbf{X})^{-}$, any generalized inverse of $(\mathbf{X}'\mathbf{X})$, and $\hat{\sigma}^2 = (\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})'(\mathbf{y} - \mathbf{X}\tilde{\boldsymbol{\beta}})/(N - r)$.

Proof. Left as an exercise. One possible approach centers on finding a linearly equivalent FR model.

Corollary 15.6.2 For $\text{GGLM}_{N,q}\text{LTFR}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta} | \mathbf{R}\boldsymbol{\beta} = \mathbf{r}, \sigma^2 \mathbf{D})$ with Gaussian errors, $\mathbf{D} = \mathbf{D}'$ ($N \times N$) is known and positive definite, $\mathbf{R}\boldsymbol{\beta} = \mathbf{r}$ is a consistent system of r equations, with $\text{rank}(\mathbf{R}) = r \leq q$. If $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is an estimable parameter and $H_0 = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$ is testable, then the LRT consists of rejecting H_0 for improbably large values of the F statistic

 (to be completed by the reader).

Proof. Left as an exercise.

Corollary 15.6.3 The $\text{GLM}_{N,q}(\mathbf{y}; \mathbf{X}\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors has estimable scalar secondary parameter $\theta = \mathbf{C}\boldsymbol{\beta}$ with $1 \times q$ $\mathbf{C} \neq \mathbf{0}$ and 1×1 $m = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}' \neq 0$. If $t = \hat{\theta}/(m\hat{\sigma}^2)^{1/2}$, then the following hold.

(a) The *likelihood ratio test* of size α for testing $H_0 = \mathbb{B}(\theta = \theta_0)$ versus $H_A = \mathbb{B}(\theta \neq \theta_0)$ is $\phi(\mathbf{y}) = \mathbb{B}(|t| > t_{\text{crit}})$ in which t_{crit} is the $100(1 - \alpha/2)$ percentile of the $t(N - q)$ distribution.

(b) The *likelihood ratio test* of size α for testing $H = \mathbb{B}(\theta > \theta_0)$ versus $H_A = \mathbb{B}(\theta \leq \theta_0)$ is $\phi(\mathbf{y}) = \mathbb{B}(t > t_{\text{crit}})$, in which t_{crit} is the $100(1 - \alpha)$ percentile of the $t(N - q)$ distribution.

For multiparameter hypotheses such as $H = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ we have a “generalized two-tail” alternative hypothesis $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$. For a one-parameter hypothesis $H = \mathbb{B}(\theta = \theta_0)$, either a two-tail or one-tail test may be chosen.

The notation of Boolean algebra (Definition 2.5) includes three operators, namely \vee (OR), \wedge (AND), as well as \neg (NOT). As throughout, $\mathbb{B}() = 1$ results from a “TRUE” argument, while $\mathbb{B}() = 0$ results from a “FALSE” argument.

Definition 15.3 (a) With \mathbf{a} an arbitrary real vector, $\mathbf{a} \in \mathfrak{R}^a$, and unknown parameter vector $\boldsymbol{\theta}$, the *composite hypothesis* is $H = \mathbb{B}(\boldsymbol{\theta} = \mathbf{0})$, with $H(\mathbf{a}) = \mathbb{B}(\mathbf{a}'\boldsymbol{\theta} = 0)$ a *component hypothesis*.
(b) Similarly, the *composite alternative* is $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \mathbf{0}) = \neg H$, and $H_A(\mathbf{a}) = \mathbb{B}(\mathbf{a}'\boldsymbol{\theta} \neq 0) = \neg H(\mathbf{a})$ is the *component alternative*.

Lemma 15.3 With the notation of the preceding definition,

$$H = \bigwedge_{\mathbf{a} \neq \mathbf{0}} H(\mathbf{a}) \tag{15.54}$$

$$H_A = \bigvee_{\mathbf{a} \neq \mathbf{0}} H_A(\mathbf{a}). \tag{15.55}$$

Proof. Here $\boldsymbol{\theta} = \mathbf{0}$ iff $\mathbf{a}'\boldsymbol{\theta} = 0 \forall \mathbf{a} \in \mathfrak{R}^a$. Consequently

$$\mathbb{B}(\boldsymbol{\theta} = \mathbf{0}) = \bigwedge_{\mathbf{a} \neq \mathbf{0}} \mathbb{B}(\mathbf{a}'\boldsymbol{\theta} = 0). \tag{15.56}$$

The results follow immediately. □

Definition 15.4 (a) Hypothesis H can be decomposed as an intersection of component hypotheses,

$$H = \bigwedge_{\mathbf{a} \neq \mathbf{0}} H(\mathbf{a}), \tag{15.57}$$

and $\phi(\mathbf{y}; \mathbf{a})$ is a test of component hypothesis $H(\mathbf{a})$ with *rejection region (critical region)* $RR(\mathbf{a}) = \{\mathbf{y} : \phi(\mathbf{y}; \mathbf{a}) = 1\}$ and *acceptance region* $AR(\mathbf{a}) = \{\mathbf{y} : \phi(\mathbf{y}; \mathbf{a}) = 0\}$.

(b) The *union-intersection test* of hypothesis H is the test specified by the acceptance region

$$RR = \bigcup_{\mathbf{a} \neq \mathbf{0}} RR(\mathbf{a}) \tag{15.58}$$

$$AR = \bigcap_{\mathbf{a} \neq \mathbf{0}} AR(\mathbf{a}). \tag{15.59}$$

Theorem 15.7 Union-Intersection Test (Roy, 1957) For $GLM_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with Gaussian errors, $\text{rank}(\mathbf{X}) = r \leq q$, and $a \times 1 \boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is an a priori testable secondary parameter [$\boldsymbol{\theta}$ is estimable and $\text{rank}(\mathbf{C}) = a \leq q$, \mathbf{C} and $\boldsymbol{\theta}_0$ are known constants]. For testing $H = \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\theta} \neq \boldsymbol{\theta}_0)$, the *union-*

intersection test (UIT) of size α is $\phi(\mathbf{y}) = \mathbb{B}[F(\mathbf{y}) > f_{\text{crit}}]$ in which

$$F(\mathbf{y}) = \{(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)'[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/a\}/\hat{\sigma}^2, \quad (15.60)$$

with

$$\tilde{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} \quad (15.61)$$

$$\hat{\boldsymbol{\theta}} = \mathbf{C}\tilde{\boldsymbol{\beta}} \quad (15.62)$$

$$\hat{\sigma}^2 = \mathbf{y}'[\mathbf{I} - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}/(N - r), \quad (15.63)$$

and f_{crit} is the $100(1 - \alpha)$ percentile of the $F(a, N - r, 0)$ distribution.

Proof. By the union-intersection principle, we define the estimator of the value of $H \equiv \mathbb{B}(\boldsymbol{\theta} = \boldsymbol{\theta}_0)$ as

$$\hat{H} = \bigwedge_{\mathbf{a} \neq \mathbf{0}} \hat{H}(\mathbf{a}). \quad (15.64)$$

An appropriate test procedure for the component sub-hypothesis, $H(\mathbf{a}) = \mathbb{B}(\mathbf{a}'\boldsymbol{\theta} = \mathbf{a}'\boldsymbol{\theta}_0)$, is a two-sided t test, $\phi(\mathbf{y}; \mathbf{a}) = \mathbb{B}(t^2(\mathbf{a}) > t_{\text{crit}}^2)$. Here t_{crit} is the $100(1 - \alpha/2)$ percentile of the $t(N - r, 0)$ distribution and $t_{\text{crit}}^2 = f_{\text{crit}}$ is the $100(1 - \alpha)$ percentile of the $F(1, N - r, 0)$ distribution. In turn,

$$t(\mathbf{a}) = \frac{\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0}{(\hat{\sigma}^2\mathbf{a}'\mathbf{M}\mathbf{a})^{1/2}}, \quad (15.65)$$

with $\hat{\mathbf{V}}(\hat{\boldsymbol{\theta}}) = \hat{\sigma}^2\mathbf{M}$ and $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$. Since $\boldsymbol{\theta}$ is estimable, \mathbf{M} is nonsingular. The acceptance region for the test is

$$\begin{aligned} AR(\mathbf{a}) &= \{\mathbf{y} : \phi(\mathbf{y}; \mathbf{a}) = 0\} \\ &= \{\mathbf{y} : t^2(\mathbf{a}) \leq f_{\text{crit}}\}, \end{aligned} \quad (15.66)$$

and the decision rule is

$$\hat{H}(\mathbf{a}) = \mathbb{B}[t^2(\mathbf{a}) \leq f_{\text{crit}}]. \quad (15.67)$$

The union-intersection decision rule for the composite test is

$$\begin{aligned} \hat{H} &= \bigwedge_{\mathbf{a} \neq \mathbf{0}} \hat{H}(\mathbf{a}) \\ &= \bigwedge_{\mathbf{a} \neq \mathbf{0}} \mathbb{B}[t^2(\mathbf{a}) \leq f_{\text{crit}}] \\ &= \mathbb{B}\{\sup_{\mathbf{a} \neq \mathbf{0}} [t^2(\mathbf{a})] \leq f_{\text{crit}}\}. \end{aligned} \quad (15.68)$$

Thus

$$\gamma(\mathbf{y}) = \sup_{\mathbf{a} \neq \mathbf{0}} [t^2(\mathbf{a})] \tag{15.69}$$

could be used as a UIT statistic. However, it is convenient to simplify the form as

$$\begin{aligned} \gamma(\mathbf{y}) &= \sup_{\mathbf{a} \neq \mathbf{0}} \left[\frac{(\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0)^2}{\hat{\sigma}^2 \mathbf{a}'\mathbf{M}\mathbf{a}} \right] \\ &= \sup_{\mathbf{a} \neq \mathbf{0}} \left[(\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0)(\mathbf{a}'\mathbf{M}\mathbf{a})^{-1} (\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0) / \hat{\sigma}^2 \right] \\ &= \frac{1}{\hat{\sigma}^2} \sup_{\mathbf{a} \neq \mathbf{0}} \left[(\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0)(\mathbf{a}'\mathbf{M}\mathbf{a})^{-1} (\mathbf{a}'\hat{\boldsymbol{\theta}} - \mathbf{a}'\boldsymbol{\theta}_0) \right] \\ &= (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1} (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) / \hat{\sigma}^2 \\ &= (\mathbf{C}\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\hat{\boldsymbol{\beta}} - \boldsymbol{\theta}_0) / \hat{\sigma}^2. \end{aligned} \tag{15.70}$$

The second from the last step uses a matrix theory result about quadratic forms (Schott, 2005, problem 9.42). We know $F(\mathbf{y}) = \gamma(\mathbf{y})/a \sim F(a, N - r, \omega)$, with

$$\omega = (\mathbf{C}\boldsymbol{\beta} - \boldsymbol{\theta}_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\boldsymbol{\beta} - \boldsymbol{\theta}_0) / \sigma^2. \tag{15.71}$$

Thus $F(\mathbf{y})$ can be used as the UIT statistic.

In summary, the union-intersection decision rule for the composite test is

$$\begin{aligned} \hat{H} &= \bigwedge_{\mathbf{a} \neq \mathbf{0}} \hat{H}(\mathbf{a}) \\ &= \mathbb{B} \left\{ \sup_{\mathbf{a} \neq \mathbf{0}} [t^2(\mathbf{a})] \leq f_{\text{crit}} \right\} \\ &= \mathbb{B} \left[F(\mathbf{y}) \leq \frac{1}{a} f_{\text{crit}} \right]. \end{aligned} \tag{15.72}$$

However, one problem remains, because the size of the test is too large. In fact,

$$\begin{aligned} \Pr\{\text{type I error}\} &= \Pr\{\hat{H} = \text{False} \mid H = \text{True}\} \\ &= \Pr\{F(\mathbf{y}) > a^{-1} f_{\text{crit}}(1, N - r)\} \\ &> \Pr\{F(\mathbf{y}) > f_{\text{crit}}(a, N - r)\} = \alpha \end{aligned} \tag{15.73}$$

because $a^{-1} f_{\text{crit}}(1, N - r) < f_{\text{crit}}(a, N - r)$. Since the individual tests for the component subhypotheses can be of any size, we need only specify that the critical value to be used throughout the proof should be $a f_{\text{crit}}(a, N - r)$ rather than $f_{\text{crit}}(1, N - r)$. □

15.4 RELATED DISTRIBUTIONS

With Gaussian errors, all components of test statistics for the univariate linear model are quadratic forms in Gaussian variables that reduce to scaled chi squares. Here we give explicit forms for the distributions.

The formulation of the likelihood ratio test makes it clear that testing a general linear hypothesis *always* corresponds to comparing two models, with the smaller nested inside the larger. The simplicity of univariate theory allows computing tests in terms of results from fitting a single model. Hence the following theorem applies to a wide variety of tests. Muller and Fetterman (2002, Chapter 5) included explicit expressions for sums of squares components of added-in-order tests and added-last tests, among others.

Theorem 15.8 With Gaussian errors and $\text{rank}(\mathbf{X}) = r \leq q$, a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ has

$$\begin{aligned} SSE/\sigma^2 &= \hat{\sigma}^2(N-r)/\sigma^2 \\ &= \mathbf{y}'[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}/\sigma^2 \\ &\sim \chi^2(N-r). \end{aligned} \quad (15.74)$$

Estimable $\boldsymbol{\theta}$ and full rank $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ ensure

$$\begin{aligned} SSH/\sigma^2 &= (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)'\mathbf{M}^{-1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/\sigma^2 \\ &\sim \chi^2[a, (\boldsymbol{\theta} - \boldsymbol{\theta}_0)'\mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)/\sigma^2]. \end{aligned} \quad (15.75)$$

Proof. Corollary 1.15 gives that $[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] = \mathbf{L}_0\mathbf{L}'_0$ is idempotent of rank $N-r$ with $\mathbf{L}'_0\mathbf{L}_0 = \mathbf{I}_{N-r}$ and $\mathbf{L}'_0\mathbf{X} = \mathbf{0}$. With $\mathbf{y} \sim \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \sigma^2\mathbf{I}_N)$ and $\boldsymbol{\beta}'\mathbf{X}'\mathbf{L}_0\mathbf{L}'_0\mathbf{X}\boldsymbol{\beta} = 0$, Theorem 9.4 gives $SSE/\sigma^2 \sim \chi^2(N-r, 0)$. If

$$\mathbf{y}_t = \mathbf{y} - \mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\boldsymbol{\theta}_0, \quad (15.76)$$

having $\mathbf{y} \sim \mathcal{N}_N(\mathbf{X}\boldsymbol{\beta}, \mathbf{I}_N\sigma^2)$ implies $\mathbf{y}_t \sim \mathcal{N}_N[\mathbf{E}(\mathbf{y}_t), \mathbf{I}_N\sigma^2]$, with

$$\mathbf{E}(\mathbf{y}_t) = \mathbf{X}\boldsymbol{\beta} - \mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\boldsymbol{\theta}_0. \quad (15.77)$$

If $\mathbf{T} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ then $\hat{\boldsymbol{\theta}} = \mathbf{T}\mathbf{y}$. Testable $\boldsymbol{\theta}$ ensures $\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X}) = \mathbf{C}$. Hence $\mathbf{T}\mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1} = \mathbf{I}_a$ and

$$\begin{aligned} \hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0 &= \mathbf{T}\mathbf{y} - \mathbf{T}[\mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}]\boldsymbol{\theta}_0 \\ &= \mathbf{T}\mathbf{y}_t. \end{aligned} \quad (15.78)$$

In turn

$$\begin{aligned} \hat{\delta} &= (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)'\mathbf{M}^{-1}(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \\ &= \mathbf{y}'_t\mathbf{T}'\mathbf{M}^{-1}\mathbf{T}\mathbf{y}_t. \end{aligned} \quad (15.79)$$

If $\mathbf{A} = \mathbf{T}'\mathbf{M}^{-1}\mathbf{T}$, then $\hat{\delta} = \mathbf{y}'_t\mathbf{A}\mathbf{y}_t$ and

$$\begin{aligned}
 A^2 &= (\mathbf{T}'\mathbf{M}^{-1}\mathbf{T})(\mathbf{T}'\mathbf{M}^{-1}\mathbf{T}) \\
 &= \mathbf{T}'\mathbf{M}^{-1}\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\{\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{C}'\}\mathbf{M}^{-1}\mathbf{T} \\
 &= \mathbf{T}'\mathbf{M}^{-1}\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\{\mathbf{C}'\}\mathbf{M}^{-1}\mathbf{T} \\
 &= \mathbf{A}.
 \end{aligned}
 \tag{15.80}$$

Furthermore, with $\boldsymbol{\theta}_{c_0} = \mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\boldsymbol{\theta}_0$,

$$\begin{aligned}
 [\mathbf{E}(\mathbf{y}_t)]'\mathbf{A}\mathbf{E}(\mathbf{y}_t) &= [\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\theta}_{c_0}]'\mathbf{T}'\mathbf{M}^{-1}\mathbf{T}[\mathbf{X}\boldsymbol{\beta} - \mathbf{X}\boldsymbol{\theta}_{c_0}] \\
 &= [\mathbf{T}\mathbf{X}\boldsymbol{\beta} - \mathbf{T}\mathbf{X}\boldsymbol{\theta}_{c_0}]'\mathbf{M}^{-1}[\mathbf{T}\mathbf{X}\boldsymbol{\beta} - \mathbf{T}\mathbf{X}\boldsymbol{\theta}_{c_0}] \\
 &= (\boldsymbol{\theta} - \boldsymbol{\theta}_0)'\mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0).
 \end{aligned}
 \tag{15.81}$$

The idempotency of \mathbf{A} combines with $\mathbf{y}_t \sim \mathcal{N}_N[\mathbf{E}(\mathbf{y}_t), \mathbf{I}_N\sigma^2]$ and Theorem 9.4 to give $\hat{\delta}/\sigma^2 \sim \chi^2[a, (\boldsymbol{\theta} - \boldsymbol{\theta}_0)'\mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)/\sigma^2]$. □

15.5 TRANSFORMATIONS AND INVARIANCE PROPERTIES

A test which does not vary when the data have been transformed in an unimportant way will usually be preferred to any test without the same invariance property. Any one of a large number of often mutually incompatible invariance properties may seem ideal, depending upon the application. However, (1) the scale of the data (nominal, ordinal, interval, ratio), (2) the mathematical constraints imposed among parameters by the hypothesis (inequalities, orderings, differences, ratios), and (3) the statistical distribution assumed combine to greatly narrow the range of sensible choices.

Our focus centers nearly all of the time on linear models of interval-scale data with Gaussian errors. Furthermore, most tests of interest involve linear relationships among parameters. Consequently, and not surprisingly, statisticians have focused on invariance under linear transformations in such settings. Transformation of the data (or functions of the data, such as parameter estimators), contrast matrices, and parameters all hold some interest.

Definition 15.5 (a) A test ϕ has *scale invariance* iff its value does not vary with changes in the units of the original observations, i.e., $\phi(\{y_{ij}\}) = \phi(\{b_j y_{ij}\})$.

(b) A test ϕ has *location invariance* iff its value its value does not vary with changes in the origin of the original observations, i.e., $\phi(\{y_{ij}\}) = \phi(\{a_j + y_{ij}\})$.

(c) A test ϕ is *linearly invariant*, which is often abbreviated as *invariant*, iff its value is scale invariant. Many, but not all, such applications also have location invariance.

For random vector \mathbf{x} , with \mathbf{A} and \mathbf{b} conforming constants, statisticians often describe the transformation of $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ as a linear transformation. However, a

precise use of mathematical definitions of transformations would classify $\mathbf{y} = \mathbf{A}\mathbf{x}$ as a linear transformation and $\mathbf{y} = \mathbf{A}\mathbf{x} + \mathbf{b}$ as a linear transformation plus a translation.

The great majority of linear models used in practice either have $\mathbf{1}_N$ in \mathbf{X} or have $\mathbf{1}_N$ spanned by the columns of \mathbf{X} (i.e., $\mathbf{1}_N = \mathbf{X}\mathbf{t}_0$ for some $q \times 1$ constant \mathbf{t}_0). In either case the columns of \mathbf{X} span an intercept (Definition 2.7). At the same time, most hypotheses tested in practice exclude the intercept ($\mathbf{C}\mathbf{t}_0 = \mathbf{0}$). The parameter deserves the name because it equals the y -axis intercept in plotting the regression function. For the univariate GLM, the usual F test is always invariant to a scale transformation of the response and predictors. In addition, if the model spans an intercept and the hypothesis excludes the intercept, the test has location invariance for the response. All other tests lack location invariance. Section 16.8 contains details and proofs for the more general multivariate case.

Multiplying all y_i by a nonzero constant k has simple consequences. Both β and θ are multiplied by k , while SSH and SSE are multiplied by k^2 . Most importantly, the F statistic, p value, and R^2 values do not change (are invariant).

Theorem 15.9 (a) Both SSH and SSE are invariant to a square, full-rank transformation of the rows of \mathbf{C} and θ_0 , with $\mathbf{C}_T = \mathbf{T}\mathbf{C}$ and $\theta_{0T} = \mathbf{T}\theta_0$.

(a) The F statistic, p value, and R^2 are invariant to the same transformation.

Proof. SSE is not a function of \mathbf{C} . If \mathbf{T} ($a \times a$) is full rank, with $\mathbf{C}_T = \mathbf{T}\mathbf{C}$,

$$\begin{aligned} SSH &= (\hat{\theta} - \theta_0)' \mathbf{M}^{-1} (\hat{\theta} - \theta_0) \\ &= (\mathbf{C}\tilde{\beta} - \theta_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\tilde{\beta} - \theta_0) \\ &= (\mathbf{C}\tilde{\beta} - \theta_0)' (\mathbf{T}^{-1}\mathbf{T})' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} \mathbf{T}^{-1}\mathbf{T} (\mathbf{C}\tilde{\beta} - \theta_0) \\ &= (\mathbf{T}'\mathbf{C}\tilde{\beta} - \mathbf{T}\theta_0)' [\mathbf{T}\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'\mathbf{T}']^{-1} (\mathbf{T}\mathbf{C}\tilde{\beta} - \mathbf{T}\theta_0) \\ &= (\mathbf{C}_T\tilde{\beta} - \theta_{0T})' [\mathbf{C}_T(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'_T]^{-1} (\mathbf{C}_T\tilde{\beta} - \theta_{0T}), \end{aligned} \quad (15.82)$$

which suffices to prove the invariance of SSH . □

Theorem 15.10 If the model spans an intercept ($\mathbf{X}\mathbf{t}_0 = \mathbf{1}_N$ for constant \mathbf{t}_0) and the hypothesis excludes the intercept ($\mathbf{C}\mathbf{t}_0 = \mathbf{0}$), the following hold.

(a) Both SSH and SSE are invariant to a location shift of the form $\mathbf{y} + \mathbf{1}_N t_1$ for constant t_1 .

(b) The F statistic, p value, and R^2 are similarly invariant.

Proof. The proof of Theorem 16.13 provides the result.

15.6 CONFIDENCE REGIONS FOR θ

Confidence regions (defined in Section 2.10) can be obtained by inverting hypothesis tests, and a confidence region can be inverted to yield a hypothesis test.

We now prove the main results which describe how to create confidence intervals. We also prove that confidence regions exist only for parameters that are testable.

Theorem 15.11 If for any $\theta_0 \in S$ there exists a size- α test, $\phi(\mathbf{y})$, of hypothesis $H(\theta_0) = \mathbb{B}(\theta = \theta_0)$, then there exists a corresponding confidence region for θ with confidence coefficient $c(\alpha) = 1 - \alpha$. Furthermore, if the acceptance region of $\phi(\mathbf{y})$ is $AR(\theta_0) = [\mathbf{y}_* : \phi(\mathbf{y}_*) = 0]$, then $R(\mathbf{y}) = [\theta_0 : \mathbf{y} \in AR(\theta_0)]$ is the corresponding confidence region.

Proof. Acceptance region $AR(\theta_0)$ is a fixed, constant region within the sample space with boundaries defined by the choice of θ_0 and α . The proof strategy is to define $R(\mathbf{y}) = [\theta_0 : \mathbf{y} \in AR(\theta_0)]$ and then show that $R(\mathbf{y})$ is a confidence region. The first step is to note that $\theta_0 \in R(\mathbf{y})$ if and only if $\mathbf{y} \in AR(\theta_0)$. It follows that $\Pr\{\theta_0 \in R(\mathbf{y})\} = \Pr\{\mathbf{y} \in AR(\theta_0)\}$. Evaluating the probability at $\theta_0 = \theta$ reveals $\Pr\{\theta \in R(\mathbf{y})\} = \Pr\{\mathbf{y} \in AR(\theta_0) | \theta_0 = \theta\} = 1 - \Pr\{\text{type I error}\} = 1 - \alpha$. \square

Theorem 15.12 If an exact $100(1 - \alpha)$ percent confidence region exists for θ , then (a) a corresponding test procedure of size α for testing $H(\theta_0) = \mathbb{B}(\theta = \theta_0)$ exists for any $\theta_0 \in S$ and (b) $H(\theta_0)$ is a *testable* hypothesis. (c) If the confidence region is $R(\mathbf{y})$, then the acceptance region of the corresponding test is $AR(\theta_0) = [\mathbf{y}_* : \theta_0 \in R(\mathbf{y}_*)]$.

Proof. For $R(\mathbf{y})$ a confidence region for θ with confidence coefficient $1 - \alpha$, $\Pr\{\theta \in R(\mathbf{y})\} = 1 - \alpha$. A particular fixed value $\theta_0 \in S$ gives $AR(\theta_0) = [\mathbf{y}_* : \theta_0 \in R(\mathbf{y}_*)]$, the set of all values in the sample space for which the confidence region would contain θ_0 . Having $AR(\theta_0)$ defined allows defining a Boolean function $\phi(\mathbf{y})$ and proving that $\phi(\mathbf{y})$ is a size- α test. For $RR(\theta_0)$ the complement of $AR(\theta_0)$, AR and RR define a mutually exclusive and together exhaustive partition of the sample space. The Boolean function $\phi(\mathbf{y})$ indicates whether y is in RR , with $\phi(\mathbf{y}) = 0$ if $y \in AR(\theta_0)$ and $\phi(\mathbf{y}) = 1$ if $y \in RR(\theta_0)$. Showing that the test is size α begins by noting $\mathbf{y} \in AR(\theta_0)$ if and only if $\theta_0 \in R(\mathbf{y})$. Now $\Pr\{\mathbf{y} \in AR(\theta_0) | \theta_0 = \theta\} = \Pr\{\theta \in R(\mathbf{y})\} = 1 - \alpha$. It follows that $\phi(\mathbf{y})$ is a size- α test of hypothesis $H(\theta_0) = \mathbb{B}(\theta = \theta_0)$. Tests do not exist for nontestable hypotheses. Therefore we conclude that $H(\theta_0)$ is testable. \square

The last two theorems and proofs reveal the relationship between tests and confidence regions. The confidence region $R(\mathbf{y}) = [\theta_0 : \mathbf{y} \in AR(\theta_0)]$ partitions the parameter space as a function of the data. On the other hand, the acceptance region $AR(\theta_0) = [\mathbf{y}_* : \theta_0 \in R(\mathbf{y}_*)]$ partitions the data space as a function of the parameter. For a given value of \mathbf{y} , confidence region $R(\mathbf{y})$ is the set of all choices of $\theta_0 \in S$ such that hypothesis $H(\theta_0)$ is not rejected. The fact that $R(\mathbf{y})$ depends on the data makes it stochastic with boundaries that depend on the realization of \mathbf{y} .

Example 15.3 In the context of a $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with Gaussian errors, inverting a test to obtain a confidence region is easy to illustrate when θ ($a \times 1$) is a scalar ($a = 1$). The two-sided t test, $\phi(\mathbf{y}) = \mathbb{B}(t^2 > f_{crit})$, of $H(\theta_0) = \mathbb{B}(\theta = \theta_0)$ has scalar $m = \mathbf{C}'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ (because $a = 1$) and acceptance region

$$\begin{aligned} AR(\theta_0) &= \left\{ \mathbf{y} : t^2(\mathbf{y}) = (\hat{\theta} - \theta_0)^2 / (\hat{\sigma}^2 m) \leq f_{crit} \right\} \\ &= \left\{ \mathbf{y} : \hat{\theta} - (f_{crit} \hat{\sigma}^2 m)^{1/2} \leq \theta_0 \leq \hat{\theta} + (f_{crit} \hat{\sigma}^2 m)^{1/2} \right\} \\ &= \left\{ \mathbf{y} : \theta_0 \in \left[\hat{\theta} \pm (f_{crit} \hat{\sigma}^2 m)^{1/2} \right] \right\}. \end{aligned} \tag{15.83}$$

It follows that the confidence region is interval

$$R(\mathbf{y}) = \left[\hat{\theta} \pm (f_{crit} \hat{\sigma}^2 m)^{1/2} \right]. \tag{15.84}$$

Here, f_{crit} is the appropriate critical value defined as the $100(1 - \alpha)$ percentile of the $F(1, N - r, 0)$ distribution and $\hat{V}(\hat{\theta}) = \hat{\sigma}^2 m$. By Theorem 15.11, the confidence coefficient for $R(\mathbf{y})$ is $1 - \alpha$.

Example 15.4 In the context of a $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with Gaussian errors, more generality arises by inverting an F test to obtain a confidence region for θ ($a \times 1$). The test $\phi(\mathbf{y}) = \mathbb{B}(F > f_{crit})$ of $H(\theta_0) = \mathbb{B}(\theta = \theta_0)$ has acceptance region

$$AR(\theta_0) = \left[\mathbf{y} : F(\mathbf{y}) = (\hat{\theta} - \theta_0)' m^{-1} (\hat{\theta} - \theta_0) / (a \hat{\sigma}^2) \leq f_{crit} \right]. \tag{15.85}$$

It follows that the $100(1 - \alpha)$ percent ellipsoidal confidence region is

$$R(\mathbf{y}) = \left[\theta_0 : (\hat{\theta} - \theta_0)' m^{-1} (\hat{\theta} - \theta_0) \leq (a \hat{\sigma}^2) f_{crit} \right]. \tag{15.86}$$

The appropriate critical value f_{crit} is defined in Theorem 15.7.

In practice, high dimensional ($a > 2$) hyperellipsoidal confidence regions are of limited utility. To facilitate graphical representations, a practical approach is to create a set of simultaneous confidence intervals for the elements of θ . Inverting a Bonferroni test, a Tukey test, a Scheffé test, or any of the other multiple-comparison test procedures gives the desired confidence region. The intervals for the individual elements of θ collectively define a high-dimensional rectangular confidence region for θ which has confidence coefficient no more than $1 - \alpha$.

Example 15.5 A $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with Gaussian errors allows many approaches to formulating simultaneous confidence intervals for the k elements of set $S_a = [\mathbf{a}'\theta : \mathbf{a} \in \{\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_k\}]$, with k finite. For each of the k choices of \mathbf{a} , the Bonferroni test of $H(\mathbf{a}, \theta_0) = \mathbb{B}(\mathbf{a}'\theta = \mathbf{a}'\theta_0)$ has acceptance region

$$AR(\mathbf{a}, \theta_0) = \left[\mathbf{y} : F(\mathbf{y}) = [\mathbf{a}'(\hat{\theta} - \theta_0)]^2 / (\hat{\sigma}^2 \mathbf{a}' m \mathbf{a}) \leq f_{crit} \right]. \tag{15.87}$$

The critical value f_{crit} depends on k and is the $100[1 - (\alpha/k)]$ percentile of $F(1, N - r, 0)$. The corresponding Bonferroni confidence interval for $\mathbf{a}'\boldsymbol{\theta}$,

$$R(\mathbf{a}, \mathbf{y}) = \left[\mathbf{a}'\hat{\boldsymbol{\theta}} \pm (f_{\text{crit}} \hat{\sigma}^2 \mathbf{a}'\mathbf{M}\mathbf{a})^{1/2} \right], \tag{15.88}$$

has a confidence coefficient no greater than $1 - \alpha$. Similarly, the Union-Intersection test (also known as the Scheffé test) can be inverted to obtain k simultaneous confidence intervals. The Scheffé test of $H(\mathbf{a}, \boldsymbol{\theta}_0) = \mathbb{E}(\mathbf{a}'\boldsymbol{\theta} = \mathbf{a}'\boldsymbol{\theta}_0)$ has acceptance region

$$AR(\mathbf{a}, \boldsymbol{\theta}_0) = \left[\mathbf{y} : F(\mathbf{y}) = \left[\mathbf{a}'(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) \right]^2 / (a\hat{\sigma}^2 \mathbf{a}'\mathbf{M}\mathbf{a}) \leq f_{\alpha, \text{crit}} \right]. \tag{15.89}$$

Here $f_{\alpha, \text{crit}}$ is the $100(1 - \alpha)$ percentile of $F(a, N - r, 0)$, which depends only on a and not on k . The Scheffé confidence interval for $\mathbf{a}'\boldsymbol{\theta}$,

$$R(\mathbf{a}, \mathbf{y}) = \left[\mathbf{a}'\hat{\boldsymbol{\theta}} \pm (af_{\alpha, \text{crit}} \hat{\sigma}^2 \mathbf{a}'\mathbf{M}\mathbf{a})^{1/2} \right], \tag{15.90}$$

has a confidence coefficient no greater than $1 - \alpha$. Unless k is roughly on the order of 2^a , the Bonferroni intervals will be narrower than the Scheffé intervals. Comparing $a \cdot f_{\alpha, \text{crit}}$ to f_{crit} illustrates the claim. When set S_a comprises k points on a regression curve or surface, then the set of confidence intervals define a confidence band. Depending on the kind of confidence band desired, k can be very small or extremely large. Stewart (1987, 1991) discussed confidence band procedures and compared various methods with respect to their graphical advantages and disadvantages.

EXERCISES

15.1 Prove the following.

Lemma. If $\mathbf{X} = \mathbf{1}_n \otimes \mathbf{X}_{\text{Es}}$, then $(\mathbf{X}'\mathbf{X})^- = n^{-1}(\mathbf{X}'_{\text{Es}}\mathbf{X}_{\text{Es}})^-$.

15.2 Model 1 is

$$\begin{aligned} \mathbf{y} &= \left(\mathbf{1}_n \otimes \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \right) \boldsymbol{\beta}_1 + \mathbf{e} \\ &= \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{e}. \end{aligned}$$

Finding a valid and unambiguous test of the general linear hypothesis $H_0 : \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$ requires verifying that (a) $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ is estimable and (b) the hypothesis is testable. In the following, with $\boldsymbol{\theta}_j = \mathbf{C}_j\boldsymbol{\beta}$, verify (a) holds or does not hold and verify (b) holds or does not hold. It is not sufficient to merely state the correct answer. You must briefly justify each positive or negative answer. For each estimable $\boldsymbol{\theta}$, describe each element very briefly as a function of cell means.

15.2.1 $\mathbf{C}_1 = [1 \ 0 \ 0 \ 0]$

15.2.2 $\mathbf{C}_2 = [0 \ 1 \ 0 \ 0]$

$$15.2.3 \mathbf{C}_3 = \begin{bmatrix} 1 & 1 & 0 & 0 \end{bmatrix}$$

$$15.2.4 \mathbf{C}_4 = \begin{bmatrix} 3 & 1 & 1 & 1 \end{bmatrix}/3$$

$$15.2.5 \mathbf{C}_5 = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

$$15.2.6 \mathbf{C}_6 = \begin{bmatrix} 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

15.3 A complete, unbalanced, G group one-way ANOVA may be written as

$$\mathbf{y} = \begin{bmatrix} \mathbf{1}_{N_1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1}_{N_G} \end{bmatrix} \begin{bmatrix} \mu_1 \\ \vdots \\ \mu_G \end{bmatrix} + \mathbf{e}.$$

Assuming i.i.d. $e_i \sim \mathcal{N}(0, \sigma^2)$, a testable hypothesis of the form $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ with $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$ may be expressed in terms of a possibly noncentral F statistic. The total sample size is $N_+ = \sum_{g=1}^G N_g$.

15.3.1 Briefly explain what “complete” requires of a particular sample for the design being considered.

15.3.2 Briefly explain what “unbalanced” requires of a particular sample for the design being considered.

15.3.3 Consider $\mathbf{C} = [1 \ -1 \ 0 \ \cdots \ 0]$, with $\boldsymbol{\theta}_0 = \mathbf{0}$, and the corresponding F test in this particular model.

(a) What are the degrees of freedom for the test in terms of $\{N_g, G\}$?

(b) Given the complete unbalanced design, what constraints on $\{N_g, G\}$ must be present for the test to be well defined?

(c) Find a convenient form for the noncentrality parameter in terms of $\{\mu_g, N_g, \sigma^2\}$.

15.3.4 In this part, consider the special case of a completely balanced design.

(a) Simplify the expression for the design matrix given the special case.

(b) If $G = 3$ and $\mathbf{C} = [\mathbf{1}_2 \ -\mathbf{I}_2]$, with $\boldsymbol{\theta}_0 = \mathbf{0}$, find a convenient form for the noncentrality parameter in terms of $\{\mu_g, N_g, \sigma^2\}$ (with a completely balanced design).

Tests for Multivariate Linear Models

16.1 MOTIVATION

Throughout the chapter, we always assume Gaussian errors in considering a multivariate linear model, a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with $\text{rank}(\mathbf{X}) = r \leq q$. With fixed, known constants $\{\mathbf{C}, \mathbf{U}, \Theta_0\}$, a secondary parameter is $a \times b$ $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$. Most of the time, but not always, we will require $\text{rank}(\mathbf{C}) = a \leq q$ and $\text{rank}(\mathbf{U}) = b \leq p$. The associated general linear hypothesis (GLH) may be stated as

$$H_0 : \Theta = \Theta_0, \tag{16.1}$$

with corresponding alternative $H_A : \Theta \neq \Theta_0$. In terms of Boolean variables, with $\mathbb{B}(\cdot) \in \{0, 1\}$, the hypothesis may be written as $H_0 = \mathbb{B}(\Theta = \Theta_0)$ versus $H_A = \mathbb{B}(\Theta \neq \Theta_0)$. Results for the univariate GLM will always be included as the special case of $p = 1$ column in \mathbf{Y} and \mathbf{B} , with $\mathbf{U} = [1]$. The first issue addressed is testability. The above hypothesis is *testable* iff the parameter Θ is estimable and the contrast matrices are full rank. We do not have hypothesis testing procedures for nontestable hypotheses.

Another important issue is whether Θ is truly a fixed constant. If the model and its parameters are defined prior to collection of the data, then Θ is an unknown, fixed constant. However, one frequently discovers interesting aspects of the data which are directly suggested by the results at hand but not considered prior to data collection. In such cases Θ may depend in some way upon the observed \mathbf{Y} . As in the univariate model (Definition 15.1), we shall refer to parameters defined independently of \mathbf{Y} as *a priori* parameters. Others, including ones obviously suggested by the data, will be called *post hoc* parameters. As indicated by the status of the parameter, the corresponding hypothesis is either an *a priori* hypothesis or a *post hoc* hypothesis. Different statistical tests are required for the two different kinds of hypotheses.

16.2 TESTABILITY OF MULTIVARIATE HYPOTHESES

In the theoretical formulation of any hypothesis test procedure, the hypothesis is always manipulated in terms of its linearly independent (LIN) components. Therefore one requires the rows of $a \times q$ \mathbf{C} to be LIN and the columns of $p \times b$ \mathbf{U} to be LIN [$\text{rank}(\mathbf{C}) = a \leq q$ and $\text{rank}(\mathbf{U}) = b \leq p$]. For a univariate model ($p = 1$), \mathbf{U} is a scalar, $\mathbf{U} = 1$. In the following, we first address the issue of testability for such full-rank \mathbf{C} and \mathbf{U} matrices. Later we consider more general \mathbf{C} and \mathbf{U} . The formal definition of a testable hypothesis implies questions of testability arise only in LTFR models.

Definition 16.1 Under the assumptions of $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors, $\text{rank}(\mathbf{X}) = r \leq q$, and $a \times b$ $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$, the hypothesis $H_0 = \mathbb{B}(\Theta = \Theta_0)$ versus $H_A = \mathbb{B}(\Theta \neq \Theta_0)$ is testable iff (1) $\text{rank}(\mathbf{C}) = a \leq q$, (2) $\text{rank}(\mathbf{U}) = b \leq p$, and (3) Θ is estimable.

In FR models, \mathbf{B} and Θ are always estimable and corresponding tests are always testable. In LTFR models $H_0 = \mathbb{B}(\mathbf{B} = \mathbf{0})$ versus $H_A = \mathbb{B}(\mathbf{B} \neq \mathbf{0})$ is never testable, while Θ may or may not be estimable. As in the univariate case, estimability and testability require only the least squares assumptions (Homogeneity, Independence, Linearity, Existence of finite second moments) and do not require a particular distribution.

To motivate the relationship between estimable parameters and testable hypotheses, we highlight the following argument. We seek test statistics which are functions of $\hat{\Theta} = \mathbf{C}\tilde{\mathbf{B}}\mathbf{U}$, with $\tilde{\mathbf{B}} = \hat{\mathbf{B}}$ in the FR case and $\mathbf{U} = 1$ in univariate models. The quantity $\hat{\Theta}$ must be invariant to the choice of $\tilde{\mathbf{B}}$ for the test to be invariant. The test result should not depend on the choice of generalized inverse! We have previously proven $\hat{\Theta}$ is invariant iff Θ is estimable. A collection of theorems in Chapters 11 and 12 give estimability criteria. In particular, $\mathbf{C}(\mathbf{X}'\mathbf{X})^{-}(\mathbf{X}'\mathbf{X}) = \mathbf{C}$ if and only if $\mathbf{C}\mathbf{B}\mathbf{U}$ is estimable.

Lemma 16.1 A $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ leads to considering constants $a \times a$ $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}\mathbf{C}'$ and $p \times b$ \mathbf{U} .

- (a) If \mathbf{X} is full rank, then $\text{rank}(\mathbf{M}) = a$ and $\text{rank}(\mathbf{U}) = b$ provide necessary and sufficient conditions for testability of $a \times b$ $\mathbf{C}\mathbf{B}\mathbf{U}$.
- (b) If \mathbf{X} is less than full rank, then the conditions $\text{rank}(\mathbf{M}) = a$ and $\text{rank}(\mathbf{U}) = b$ provide necessary but not sufficient conditions for testability.
- (c) If \mathbf{X} is less than full rank, then the conditions $\text{rank}(\mathbf{M}) = a$ and $\text{rank}(\mathbf{U}) = b$ combined with the requirement $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-}(\mathbf{X}'\mathbf{X})$, or any other condition which guarantees estimability, provide a necessary and sufficient set of conditions for testability.

Proof. Left as an exercise.

Certain related conditions are worth remembering. With full rank \mathbf{X} , the rank of \mathbf{M} is the rank of \mathbf{C} . No matter the rank of \mathbf{X} , with \mathbf{C} $a \times q$, both $\text{rank}(\mathbf{C})$ and $\text{rank}(\mathbf{M})$ are necessarily bounded above by $\min(a, q)$.

Theorem 16.1 A nontestable hypothesis has no reasonable test. For $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and $\text{rank}(\mathbf{X}) = r \leq q$, nonestimable $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ ($a \times b$) may have $\text{rank}(\mathbf{C}) = a \leq q$ and $\text{rank}(\mathbf{U}) = b \leq p$. Even with $\text{rank}(\mathbf{C}) = a$ and $\text{rank}(\mathbf{U}) = b$ the hypothesis $H_0 = \mathbb{B}(\Theta = \Theta_0)$ versus $H_A = \mathbb{B}(\Theta \neq \Theta_0)$ is not testable. No reasonable test can be found; i.e., any test procedure will have size $\alpha = 0$ and power $1 - \beta = 0$. Furthermore the test will not be invariant to changes in irrelevant quantities.

Proof. The univariate model is a special case of the multivariate model. Hence the proof of the nonexistence of a valid test in the univariate case demonstrates a valid test can not be guaranteed in the multivariate case. \square

Corollary 16.1 If $\Theta_s = [\Theta'_1 \ \Theta'_2]'$ and Θ_1 is estimable but Θ_2 is not, then tests for Θ_s are indistinguishable from tests for Θ_1 .

Proof. The previous chapter has a proof for the univariate case. The \mathbf{U} matrix plays no role in estimability, although it does affect testability. \square

The multivariate test statistics described in Chapter 3 are all functions of

$$\hat{\Delta} = (\mathbf{C}\hat{\mathbf{B}}\mathbf{U} - \Theta_0)' [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1} (\mathbf{C}\hat{\mathbf{B}}\mathbf{U} - \Theta_0). \tag{16.2}$$

The $b \times b$ matrix generalizes $SSH = \hat{\delta}$ from the univariate setting, which leads to the alternate notation $S_h = \hat{\Delta}$. In the FR case $\hat{\mathbf{B}} = \mathbf{B}$ and $\mathbf{X}'\mathbf{X}$ is nonsingular. The test statistics depend on \mathbf{C} and Θ_0 only through $\hat{\Delta}$. In a univariate model ($p = 1$ or $b = 1$), $\hat{\Delta}$ is a scalar. The ability to compute $\hat{\Delta}$ depends on the existence of $[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']^{-1}$. The value of $\hat{\Delta}$ can be computed iff $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}' \neq \mathbf{0}$ has full rank. Testable Θ guarantees full-rank \mathbf{M} .

In many cases it is possible to compute $\hat{\Delta}$ (and therefore the test statistic) even though Θ is not estimable and the hypothesis is not testable. The result is formally stated in the following theorem. Subsequently, we prove a theorem which answers the question ‘‘What hypothesis is the test statistic addressing?’’ in such cases.

Theorem 16.2 The $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors has $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ $a \times b$, $\text{rank}(\mathbf{X}) = r \leq q$, and $\text{rank}(\mathbf{C}) = a \leq q$, and $\text{rank}(\mathbf{U}) = b \leq p$.

- (a) Θ is estimable implies $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ has full rank, although
- (b) $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$ having full rank does not imply Θ is estimable.

The preceding chapter includes a proof in the univariate case. As always, \mathbf{U} plays no role in determining estimability.

Theorem 16.3 For a $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with Gaussian errors, $\text{rank}(X) = r < q$, $\text{rank}(C) = a \leq q$, and $\text{rank}(U) = b \leq p$, $\Theta = CBU$ ($a \times b$) may not be estimable even though $M^{-1} = [C(X'X)^-C']^{-1}$ exists. If so, test statistics based on

$$\hat{\Delta} = (C\tilde{B}U - \Theta_0)' M^{-1} (C\tilde{B}U - \Theta_0) \tag{16.3}$$

can be computed and the testable hypothesis actually tested is $H_0 = \mathbb{B}(\Theta_G = 0)$ versus $H_A = \mathbb{B}(\Theta_G \neq 0)$, with $\Theta_G = C(X'X)^-(X'X)BU - \Theta_0$.

Proof. Searle (1971, p. 195) provided a proof in the univariate case.

Without estimability, Θ_G is not necessarily invariant to the choice of $(X'X)^-$. For each choice of generalized inverse a potentially different testable hypothesis is tested! An example of such a poorly defined test is given in the previous chapter.

We now consider more general definitions of $a \times q$ C and $p \times b$ U . Previously we have required $\text{rank}(C) = a \leq q$ and $\text{rank}(U) = b \leq p$. Can the requirement be relaxed? We are occasionally interested in matrices C and U which are not of full rank. In any case we will need $C = AX$ (for some A) to have any hope of formulating a testable hypothesis. Otherwise Θ will not be estimable.

Estimability of Θ does not require C to have full rank. If $a \times b$ $C_1 B U_1 = \Theta_1$ is estimable, then so is the $2a \times 2b$ parameter

$$\begin{bmatrix} C_1 \\ C_1 \end{bmatrix} B [U_1 \ U_1] = \begin{bmatrix} C_1 B \\ C_1 B \end{bmatrix} [U_1 \ U_1] \tag{16.4}$$

$$C_2 B U_2 = \begin{bmatrix} \Theta_1 & \Theta_1 \\ \Theta_1 & \Theta_1 \end{bmatrix} = \Theta_2. \tag{16.5}$$

However, Θ_2 is not testable because $\text{rank}(C_2) = \text{rank}(C_1) = a < 2a$ and $\text{rank}(U_2) = \text{rank}(U_1) = b < 2b$. Also $H_{0,2} = \mathbb{B}(\Theta_2 = 0) \Leftrightarrow \mathbb{B}[(\Theta_1 = 0) \cap (\Theta_1 = 0) \cap (\Theta_1 = 0) \cap (\Theta_1 = 0)]$ is overstated and identical to $H_{0,1} = \mathbb{B}(\Theta_1 = 0)$. We should test subhypothesis $H_{0,1}$ rather than $H_{0,2}$.

Theorem 16.4 A test with a LTFR C matrix cannot be distinguished from a test with C replaced by a full (row) rank matrix with rows that span all rows of C . More specifically, $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ has Gaussian errors, $\text{rank}(X) = r \leq q$, estimable $\Theta = CBU$ ($a \times b$), $\text{rank}(C) = a_1 < \min(a, q)$, $\text{rank}(U) = b \leq p$,

$$\begin{aligned} \Theta &= CBU \\ &= \begin{bmatrix} C_1 \\ C_2 \end{bmatrix} BU = \begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix}, \end{aligned} \tag{16.6}$$

and $\Theta_0 = [\Theta'_{0,1} \ \Theta'_{0,2}]'$. For $j \in \{1, 2\}$, C_j , Θ_j , and $\Theta_{0,j}$ have a_j rows and $\text{rank}(C) = \text{rank}(C_1) = a_1 < \min\{a, q\}$.

(a) Tests about Θ cannot be distinguished from tests about Θ_1 ; i.e., $\Theta = \Theta_0$ iff

$\Theta_1 = \Theta_{0,1}$, and $H_0 = \mathbb{B}(\Theta = \Theta_0) = \mathbb{B}(\Theta_1 = \Theta_{0,1})$.

(b) Tests about $\Theta_1 = C_1BU$ are testable; i.e., C_1 has full row rank, U has full column rank, and Θ_1 is estimable.

Proof. The rows of C_2 are linear combinations of the rows of C_1 ; i.e., $C_2 = AC_1$ for some fixed matrix A . Hence $\Theta_2 = AC_1BU = A\Theta_1$ and $\Theta = \Theta_0$ iff $\Theta_1 = \Theta_{0,1}$. Also

$$\begin{aligned} \Theta = CBU &= \begin{bmatrix} I \\ A \end{bmatrix} C_1BU \\ &= \begin{bmatrix} \Theta_1 \\ \Theta_2 \end{bmatrix} = \begin{bmatrix} I \\ A \end{bmatrix} \Theta_1. \end{aligned} \tag{16.7}$$

Estimable Θ implies Θ_1 is estimable; combining estimability with full rank of C_1 implies $\hat{\Delta}_1$ exists. While $\hat{\Delta} = (\Theta - \Theta_0)'[C(X'X)^{-1}C']^{-1}(\hat{\Theta} - \Theta_0)$ does not exist and cannot be computed, $\hat{\Delta}_1 = (\hat{\Theta}_1 - \Theta_{0,1})'[C_1(X'X)^{-1}C_1']^{-1}(\hat{\Theta}_1 - \Theta_{0,1})$ does exist and can be computed. The commonly used tests require the existence of $\hat{\Sigma}_*^{-1} = (U'\hat{\Sigma}U)^{-1}$. Full column rank of U (and the assumption of full-rank Σ) guarantees the desired property. \square

Theorem 16.5 A test with a LTFR U matrix is indistinguishable from a test with U replaced by a full (column) rank matrix with columns that span all columns of U . More specifically, a $GLM_{N,p,q}(Y_i; X_iB, \Sigma)$ has Gaussian errors, $\text{rank}(X) = r \leq q$, estimable $\Theta = CBU$ ($a \times b$), and $\text{rank}(C) = a \leq q$. For $j \in \{1, 2\}$ U_j $p \times b_j$, $U = [U_1 U_2]$, and Θ_j $a \times b_j$, $\Theta = [\Theta_1 \Theta_2]$, $\Theta_0 = [\Theta_{0,1} \Theta_{0,2}]$, and $\text{rank}(U) = \text{rank}(U_1) = b_1 \leq \min\{b, p\}$.

(a) Tests about Θ are indistinguishable from tests about Θ_1 ; i.e., $\Theta = \Theta_0$ iff $\Theta_1 = \Theta_{0,1}$, and $H_0 = \mathbb{B}(\Theta = \Theta_0) = \mathbb{B}(\Theta_1 = \Theta_{0,1})$.

(b) Tests about $\Theta_1 = CBU_1$ are testable; i.e., C has full row rank, U_1 has full column rank, and Θ_1 is estimable.

Proof. (a) All columns of U_2 are linear combinations of columns of U_1 : $U_2 = U_1A$ for fixed A . Here $\Theta = \Theta_0$ iff $\Theta_1 = \Theta_{0,1}$, because $\Theta_2 = CBU_2 = CBU_1A = \Theta_1A$. (b) Matrix $\hat{\Delta}_1 = (\hat{\Theta}_1 - \Theta_{0,1})'[C(X'X)^{-1}C']^{-1}(\hat{\Theta}_1 - \Theta_{0,1})$ exists and can be computed. The commonly used tests require the existence of $\hat{\Sigma}_*^{-1} = (U_1'\hat{\Sigma}U_1)^{-1}$, which is guaranteed by the full column rank of U_1 , even though $(U'\hat{\Sigma}U)$ is singular and $(U'\hat{\Sigma}U)^{-1}$ does not exist. \square

The following theorem justifies using a simpler model which is linearly equivalent. The approach allows avoiding some of the pitfalls with LTFR models.

Theorem 16.6 Any primary or secondary expected- value parameter testable in $GLM_{N,p,q}(Y_i; X_iB, \Sigma)$ is also testable in a linearly equivalent model.

Proof. The result follows from the parallel result about estimability.

16.3 TESTS OF A PRIORI HYPOTHESES

The present chapter focuses on the GLH, $H_0 : \Theta = \Theta_0$, for a multivariate $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with $\text{rank}(X) = r \leq q$ and Gaussian errors. Here we consider only testable $\Theta = CBU$ of dimensions $a \times b$. Testable parameters must be estimable and have $\text{rank}(C) = a \leq q$, and $\text{rank}(U) = b \leq p$, for fixed and known C, U , and Θ_0 . The primary parameters B and Σ have estimators

$$\tilde{B} = (X'X)^{-1} X'Y \tag{16.8}$$

$$\hat{\Sigma} = Y'[I - X(X'X)^{-1}X']Y / (N - r). \tag{16.9}$$

Next we derive tests of $H_0 = \mathbb{B}(\Theta = \Theta_0)$ versus $H_A = \mathbb{B}(\Theta \neq \Theta_0)$ in terms of secondary parameters

$$\Theta = CBU \tag{16.10}$$

$$\Sigma_* = U'\Sigma U, \tag{16.11}$$

which are $a \times b$ and $b \times b$ respectively. With $M = C(X'X)^{-1}C'$ ($a \times a$), fully specifying distributions under the alternative involves the $b \times b$ shift matrix

$$\Delta = (\Theta - \Theta_0)'M^{-1}(\Theta - \Theta_0) \tag{16.12}$$

or the noncentrality matrix

$$\Omega = \Delta\Sigma_*^{-1}. \tag{16.13}$$

Corresponding estimators, with $\nu_e = N - \text{rank}(X)$, are

$$\hat{\Theta} = C\tilde{B}U \tag{16.14}$$

$$\nu_e^{-1}S_e = \hat{\Sigma}_* = U'\hat{\Sigma}U \tag{16.15}$$

$$S_h = \hat{\Delta} = (\hat{\Theta} - \Theta_0)'M^{-1}(\hat{\Theta} - \Theta_0). \tag{16.16}$$

Also

$$\hat{\Theta} - \Theta_0 \sim \mathcal{N}_{a,b}(CBU - \Theta_0, M, \Sigma_*) \tag{16.17}$$

$$S_e \sim \mathcal{W}_b(\nu_e, \Sigma_*, \mathbf{0}) \tag{16.18}$$

$$S_h \sim \mathcal{W}_b(a, \Sigma_*, \Delta). \tag{16.19}$$

The parameters determining the distributions of the secondary parameter estimators make it obvious that any reasonable test statistic must be a function of $\hat{\Delta}$ and $\hat{\Sigma}_*$, which generalize SSH and $\hat{\sigma}^2$ (equivalent to SSE) and reduce to them if $p = 1$. Properties of special cases help motivate the choice of statistics. Most importantly, if $b = 1$ (always true for $p = 1$), then the $b \times b$ matrix $\hat{\Omega}/a = \hat{\Delta}\hat{\Sigma}_*^{-1}/a$ is 1×1 , a scalar, and exactly equals the usual F random variable from univariate theory. In the same case, $\Omega = \Delta\Sigma_*^{-1}$ is 1×1 and reduces to $\omega = (\theta - \theta_0)'M^{-1}(\theta - \theta_0)\sigma_*^{-2}$, the F noncentrality. More generally, if

$s = \min(a, b) = 1$, then the $b \times b$ matrix $\widehat{\Delta} \widehat{\Sigma}_*^{-1} / a$ has a single nonnegative eigenvalue which is a constant multiple of an F random variable.

The most general case, with $s = \min(a, b) > 1$, creates more complication. In contrast to the univariate results, which are a special case, the multivariate likelihood ratio and union-intersection tests differ. Furthermore, other criteria, such as the substitution principle, lead to distinct tests. A form of linear invariance and exact test size stand as the properties demanded of all candidate statistics.

16.4 LINEAR INVARIANCE

A multivariate test is said to be linearly invariant if the hypothesis test (interpreted as a decision function) does not vary under full-rank transformation of the response variables being tested. Transformations of the underlying model provide a straightforward way to formalize the idea. For testable $\Theta = \mathbf{C} \mathbf{B}_U$, a test of $H_0 : \mathbf{C} \mathbf{B}_U = \Theta$ in $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with Gaussian errors may be expressed in terms of a model of $\mathbf{Y}_U = \mathbf{Y} \mathbf{U}$ ($N \times b$) by the transformation

$$\begin{aligned} \mathbf{Y} \mathbf{U} &= \mathbf{X} \mathbf{B}_U + \mathbf{E} \mathbf{U} \\ \mathbf{Y}_U &= \mathbf{X} \mathbf{B}_U + \mathbf{E}_U. \end{aligned} \tag{16.20}$$

With $\Sigma_* = \mathbf{U}' \Sigma \mathbf{U}$, the transformation implicitly defines the $\text{GLM}_{N,b,q}(\mathbf{Y}_{U_i}; \mathbf{X}_i \mathbf{B}_U, \Sigma_*)$. In turn, testing $H_0 : \mathbf{C} \mathbf{B}_U = \Theta_0$ is obviously equivalent to testing $H_0 : \mathbf{C} \mathbf{B}_U = \Theta_0$.

Definition 16.2 A $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with Gaussian errors and testable $\Theta = \mathbf{C} \mathbf{B}_U$ allows testing $H_0 : \mathbf{C} \mathbf{B}_U = \Theta_0$.

(a) If $\mathbf{Y}_U = \mathbf{Y} \mathbf{U}$ and $\mathbf{B}_U = \mathbf{B} \mathbf{U}$, the test is *linearly invariant* whenever the same test occurs for $H_0 : \Theta \mathbf{T} = \Theta_0 \mathbf{T}$ in $\text{GLM}_{N,b,q}(\mathbf{Y}_{U_i} \mathbf{T}; \mathbf{X}_i \mathbf{B}_U, \mathbf{T}' \Sigma_* \mathbf{T})$, with constant and full-rank \mathbf{T} ($b \times b$). As always, $\Sigma_* = \mathbf{U}' \Sigma_* \mathbf{U}$.

(b) The test has *location invariance* whenever the model spans an intercept ($\mathbf{X} \mathbf{t}_0 = \mathbf{1}_N$) and the hypothesis excludes it ($\mathbf{C} \mathbf{t}_0 = \mathbf{0}$).

Location invariance also can be expressed in terms of the transformation $\mathbf{Y}_{U_i} + \mathbf{1}_N \mathbf{t}'_1$ for constant \mathbf{t}'_1 ($b \times 1$). Clearly $E(\mathbf{Y}_{U_i} + \mathbf{1}_N \mathbf{t}'_1) = (\mathbf{X}_i \mathbf{B}_U + \mathbf{1}_N \mathbf{t}'_1)$. However, the specific impact on $\mathbf{B}_U = \mathbf{B} \mathbf{U}$ varies with the design. In the special case of a full-rank model with $\mathbf{1}_N$ as the first column (defining the intercept parameter), \mathbf{t}'_1 is added to the first row of \mathbf{B}_U .

Some authors include location invariance as part of linear invariance. In any case, the only functions of $\widehat{\Delta}$ and $\widehat{\Sigma}_*$ which allow achieving the invariance properties are the eigenvalues of $\widehat{\Omega} = \widehat{\Delta} \widehat{\Sigma}_*^{-1}$.

Although many appealing size- α tests exist, for general Σ , no uniformly most powerful (UMP) test of size α (among scale-invariant tests) can be found for the general linear multivariate hypothesis. Next we consider four commonly used test procedures all of which are invariant, unbiased, and admissible.

16.5 FOUR MULTIVARIATE TEST STATISTICS

The multivariate (MULTIREP) statistics involve $\nu_e = N - \text{rank}(\mathbf{X})$, \mathbf{M} ($a \times a$) with $\text{rank}(\mathbf{M}) = a$, $\widehat{\Delta}$ ($b \times b$) with $\text{rank}(\widehat{\Delta}) = s = \min\{a, b\}$, and $\text{rank}(\widehat{\Sigma}_*) = b$. The four multivariate test statistics of interest are defined in terms of the $s = \min\{a, b\}$ nonzero eigenvalues of $\widehat{\Delta}\widehat{\Sigma}_*^{-1}$, which is usually not symmetric. The symmetric and full-rank nature of $\widehat{\Sigma}_*$ ensures full-rank \mathbf{F} exists such that $\widehat{\Sigma}_* = \mathbf{F}\mathbf{F}'$ and $\widehat{\Sigma}_* = \mathbf{F}^{-t}\mathbf{F}^{-1}$. The eigenvalues of $\widehat{\Delta}\widehat{\Sigma}_*^{-1}$ coincide with the eigenvalues of the symmetric matrix $\mathbf{F}^{-1}\widehat{\Delta}\mathbf{F}^{-t}$, while the eigenvectors differ by the transformation $\mathbf{V}_F = \mathbf{F}^{-1}\mathbf{V}$. The coincidence may be proven by observing

$$\begin{aligned} (\widehat{\Delta}\widehat{\Sigma}_*^{-1} - \lambda\mathbf{I}_b)\mathbf{v} &= \mathbf{0} \\ \mathbf{F}^{-1}(\widehat{\Delta}\mathbf{F}^{-t}\mathbf{F}^{-1} - \lambda\mathbf{I}_b)(\mathbf{F}\mathbf{F}^{-1})\mathbf{v} &= \mathbf{F}^{-1}\mathbf{0} \\ (\mathbf{F}^{-1}\widehat{\Delta}\mathbf{F}^{-t} - \lambda\mathbf{I}_b)(\mathbf{F}^{-1}\mathbf{v}) &= \mathbf{0}. \end{aligned} \tag{16.21}$$

With probability 1, $\widehat{\Sigma}_*^{-1}$ is positive definite and $\widehat{\Delta}$ is either positive definite or positive semidefinite, and both are symmetric. Hence the eigenvalues in the last equation are real and nonnegative ($\mathbf{F}^{-1}\widehat{\Delta}\mathbf{F}^{-t}$ and $\widehat{\Delta}$ are congruent and therefore have the same number of positive, negative, and zero eigenvalues). The $b \times b$ matrices $\widehat{\Delta}\widehat{\Sigma}_*^{-1}$ and $\mathbf{F}^{-1}\widehat{\Delta}\mathbf{F}^{-t}$ have rank $s = \min\{a, b\}$ because $\widehat{\Sigma}_*$ and \mathbf{F} are rank b , while $\text{rank}(\widehat{\Delta}) = a = \text{rank}(\mathbf{M})$. Various authors discuss matrices which differ by a simple multiple (Kuhfeld, 1986). The matrix $\mathbf{S}_e = \nu_e\widehat{\Sigma}_*$ ($b \times b$) is the sum of squares of error and reduces to *SSE* in the univariate case. Considering $\mathbf{S}_h\mathbf{S}_e^{-1} = \widehat{\Delta}\widehat{\Sigma}_*^{-1}/\nu_e$, the multivariate generalization of *SSH/SSE*, is often convenient. The multivariate analog of the total sum of squares is $\mathbf{S}_t = \mathbf{S}_h + \mathbf{S}_e$.

The test statistics can also be defined in terms of the eigenvalues of $\mathbf{S}_h\mathbf{S}_t^{-1}$, which reduces to *SSH/SST*, or $\mathbf{S}_e\mathbf{S}_t^{-1}$, which reduces to *SSE/SST*. Writing $\mathbf{S}_h = \widehat{\Delta}$ leads to Table 16.1, which summarizes algebraic relationships among the different sets of eigenvalues. Table 1 in Muller, LaVange, Ramey, and Ramey (1992) contains additional information. Here $0 \leq \widehat{\rho}_k \leq 1$ is one of s generalized canonical correlations. It reduces to the sole multiple correlation for a univariate model which spans an intercept combined with a hypothesis excluding the intercept.

Table 16.1 Eigenvalues of Matrix Labeling Column as a Function of Eigenvalues of Matrix Labeling Row

	$\mathbf{S}_h\mathbf{S}_e^{-1}$	$\mathbf{S}_h\mathbf{S}_t^{-1}$	$\mathbf{S}_e\mathbf{S}_t^{-1}$
$\mathbf{S}_h\mathbf{S}_e^{-1}$	$\widehat{\phi}_k$	$\widehat{\phi}_k/(1 + \widehat{\phi}_k)$	$1/(1 + \widehat{\phi}_k)$
$\mathbf{S}_h\mathbf{S}_t^{-1}$	$\widehat{\rho}_k^2/(1 - \widehat{\rho}_k^2)$	$\widehat{\rho}_k^2$	$(1 - \widehat{\rho}_k^2)$
$\mathbf{S}_e\mathbf{S}_t^{-1}$	$(1 - \widehat{\lambda}_k)/\widehat{\lambda}_k$	$(1 - \widehat{\lambda}_k)$	$\widehat{\lambda}_k$

Definition 16.3 Wilks lambda ($\widehat{\Lambda}$), the likelihood ratio statistic, is

$$\begin{aligned} \widehat{\Lambda} &= |\mathbf{S}_h \mathbf{S}_e^{-1} + \mathbf{I}_b|^{-1} \\ &= |\mathbf{S}_e| / |\mathbf{S}_h + \mathbf{S}_e| = |\mathbf{S}_e| / |\mathbf{S}_t| \\ &= |\mathbf{S}_e \mathbf{S}_t^{-1}| \\ &= \prod_{k=1}^b (1 - \widehat{\rho}_k^2). \end{aligned} \tag{16.22}$$

Definition 16.4 Roy's largest root (RLR), the union-intersection principle statistic, is $\text{RLR} = \widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1})$, which is the largest eigenvalue of $\mathbf{S}_h \mathbf{S}_e^{-1}$, say $\widehat{\phi}_1$. The variable $\widehat{\phi}_1$ is a one-to-one function of and hence equivalent to $\widehat{\rho}_1^2$.

Definition 16.5 The Hotelling-Lawley trace (HLT, ANOVA analog) is

$$\text{HLT} = \text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1}) = \sum_{k=1}^b \widehat{\phi}_k = \sum_{k=1}^b \widehat{\rho}_k^2 / (1 - \widehat{\rho}_k^2). \tag{16.23}$$

Definition 16.6 The Pillai-Bartlett trace (PBT), the R^2 analog statistic, is

$$\text{PBT} = \text{tr}(\mathbf{S}_h \mathbf{S}_t^{-1}) = \sum_{k=1}^b \widehat{\rho}_k^2. \tag{16.24}$$

Pillai (1955) justified the PBT on heuristic grounds. It is linearly invariant, and the distribution depends only on the dimensions of the problem under the null. The other three statistics considered in the chapter share the same properties.

Roy (1957) proposed $\widehat{\phi}_1$ as a multivariate test statistic and derived the central distribution of $\widehat{\rho}_1^2 = \widehat{\phi}_1 / (1 + \widehat{\phi}_1)$. Some authors (including Muller and Peterson, 1984) refer to $\widehat{\rho}_1^2$ as Roy's largest root.

Lawley (1938) and Hotelling (1951) proposed T^2 as a test statistic and derived the central distribution of $T^2 / \nu_e = \text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1})$. Some authors refer to $\text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1})$ as "the trace criterion" or as "Hotelling's trace criterion" (Timm, 1975).

Timm (1975, p. 148–149) noted the noncentral distributions have been very difficult to compute, although recent advances in algorithms have led to published tables. Muller and Peterson (1984) also discussed the noncentral distributions.

For a multivariate model, $\mathbf{S}_e \sim \mathcal{W}_b(\nu_e, \boldsymbol{\Sigma}_*)$, independently of

$$\mathbf{S}_h = \widehat{\Delta} \sim \mathcal{W}_b(a, \boldsymbol{\Sigma}_*, \Delta). \tag{16.25}$$

Under H_0 , $\Delta = \mathbf{0}$ and $\widehat{\Delta}$ has a central Wishart distribution. Aside from certain special cases, addressed later, the distribution of $\widehat{\Lambda}$ is complicated.

Theorem 16.7 Under H_0 , the statistic $\widehat{\Lambda}$ equals a product of b independent beta-1 random variables. More specifically, if $\mathbf{S}_e \sim \mathcal{W}_b(\nu_e, \boldsymbol{\Sigma})$ independently of $\mathbf{S}_h \sim \mathcal{W}_b(a, \boldsymbol{\Sigma})$ and $\nu_e \geq b$, then

$$\widehat{\Lambda} = |\mathbf{S}_h \mathbf{S}_e^{-1} + \mathbf{I}|^{-1} = \prod_{k=1}^b x_k, \tag{16.26}$$

with $\{x_k\}$ independent beta random variables, $x_k \sim \beta[(\nu_e+k-b)/2, a/2]$. If $a = 1$, $\widehat{\Lambda}$ has the same distribution as a single beta random variable with parameters $\{(\nu_e + 1 - b)/2, b/2\}$ and hence is a one-to-one function of an F .

Proof. Rao (1973, Section 8.b.2.xi, Theorem 3.4.3, Corollary 3.4.3) gave a proof.

Corollary 16.7 With $s = \min\{a, b\}$,

(a) if $s = 1$, then

$$F = (\widehat{\Lambda}^{-1} - 1)(\nu_e + 1 - b)/(|a - b| + 1) \sim F(|a - b| + 1, \nu_e + 1 - b), \tag{16.27}$$

(b) if $s = 2$, then

$$F = (\widehat{\Lambda}^{-1} - 1)(\nu_e + 1 - b)/(|a - b| + 2) \sim F[2|a - b| + 8, 2(\nu_e + 1 - b)]. \tag{16.28}$$

Proof. Morrison (1990) provided a proof.

The results can be used for an exact LRT when $s = 1$ or $s = 2$, which covers a useful proportion of hypotheses tested in practice.

Theorem 16.8 The likelihood ratio test (LRT) statistic is $\widehat{\Lambda}$. More specifically, $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors has estimable $\boldsymbol{\Theta} = \mathbf{C} \mathbf{B} \mathbf{U}$ $a \times b$, $\text{rank}(\mathbf{X}) = r \leq q$, $\text{rank}(\mathbf{C}) = a \leq q$, and $\text{rank}(\mathbf{U}) = b \leq p$, \mathbf{C} , \mathbf{U} , and $\boldsymbol{\Theta}_0$ known constants. The LRT of $H_0 = \mathbb{B}(\boldsymbol{\Theta} = \boldsymbol{\Theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\Theta} \neq \boldsymbol{\Theta}_0)$ is $\widehat{H}_A = \mathbb{B}(\widehat{\Lambda} < t_0)$ with $\widehat{\Lambda} = |\mathbf{S}_h \mathbf{S}_e^{-1} + \mathbf{I}|^{-1} = |\mathbf{S}_e|/|\mathbf{S}_h + \mathbf{S}_e|$ and appropriate critical value t_0 . Here $\widehat{\boldsymbol{\Theta}} = \mathbf{C} \widehat{\mathbf{B}} \mathbf{U}$, $\mathbf{S}_h = (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)' [\mathbf{C}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}]^{-1} (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)$, $\widehat{\boldsymbol{\Sigma}} = \mathbf{Y}' [\mathbf{I} - \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'] \mathbf{Y} / (N - r)$, and $\mathbf{S}_e = \mathbf{U}' \widehat{\boldsymbol{\Sigma}} \mathbf{U} / (N - r)$.

Proof. Anderson (2004) included a proof. The result follows from expressing the test in terms of the unconstrained and constrained likelihood expressions in the proofs for MLEs of \mathbf{B} and $\boldsymbol{\Theta}$. The proof directly generalizes the univariate proof.

Hypothesis H_0 is rejected for improbably small values of $\widehat{\Lambda}$, relative to an appropriate critical value.

Theorem 16.9 The union-intersection (UI) test statistic is $\widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1})$. The $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors has estimable $\boldsymbol{\Theta} = \mathbf{C} \mathbf{B} \mathbf{U}$ $a \times b$, $\text{rank}(\mathbf{X}) = r \leq q$, $\text{rank}(\mathbf{C}) = a \leq q$, and $\text{rank}(\mathbf{U}) = b \leq p$, with \mathbf{C} , \mathbf{U} , and $\boldsymbol{\Theta}_0$ known constants. For testing $H_0 = \mathbb{B}(\boldsymbol{\Theta} = \boldsymbol{\Theta}_0)$ versus $H_A = \mathbb{B}(\boldsymbol{\Theta} \neq \boldsymbol{\Theta}_0)$ the UI test of H_0 versus H_A is $\widehat{H}_A = \mathbb{B}[\widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1}) > l_\alpha]$ and has test statistic $\widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1})$, the largest eigenvalue of $\mathbf{S}_h \mathbf{S}_e^{-1}$ and sometimes described as Roy's largest root statistic. The appropriate critical value is l_α . Hypothesis H_0 is rejected for improbably large values of the test statistic relative to l_α , the appropriate critical value.

Proof. Any hypothesis involving $\boldsymbol{\Theta}_0 \neq \mathbf{0}$ can be converted to a test with $\boldsymbol{\Theta}_0 = \mathbf{0}$. The proof (in the section on invariance below) is based on the fact $\boldsymbol{\Theta} = \mathbf{0}$ iff $\boldsymbol{\Theta} \mathbf{b} = \mathbf{0} \forall \mathbf{b} \in \mathfrak{R}^b$ (\mathbf{b} is $b \times 1$). In terms of Boolean algebra we have the following decomposition of the null hypothesis into infinitely many univariate hypotheses:

$$\begin{aligned} H_0 &= \mathbb{B}(\boldsymbol{\Theta} = \mathbf{0}) \\ &= \bigwedge_{\mathbf{b} \neq \mathbf{0}} \mathbb{B}(\boldsymbol{\Theta} \mathbf{b} = \mathbf{0}) \\ &= \bigwedge_{\mathbf{b} \neq \mathbf{0}} H_0(\mathbf{b}). \end{aligned} \tag{16.29}$$

For a given $\mathbf{b} \neq \mathbf{0}$, $H_0(\mathbf{b})$ is a univariate joint hypothesis and $\boldsymbol{\Theta} \mathbf{b}$ is an $(a \times 1)$ vector with $\widehat{\boldsymbol{\Theta} \mathbf{b}} \sim \mathcal{N}_a[\boldsymbol{\Theta} \mathbf{b}, \mathbf{b}' \boldsymbol{\Sigma} \mathbf{b} \mathbf{C} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}']$. The result is easily derived from $\mathbf{Y} \mathbf{b} \sim \mathcal{N}_N(\mathbf{X} \mathbf{B} \mathbf{b}, \mathbf{b}' \boldsymbol{\Sigma} \mathbf{b} \mathbf{I}_N)$. The univariate UI test of $H_0(\mathbf{b}) \equiv \mathbb{B}(\boldsymbol{\Theta} \mathbf{b} = \mathbf{0})$ is based upon the statistic

$$\begin{aligned} F(\mathbf{b}; \mathbf{Y}) &= \frac{\mathbf{b}' \mathbf{S}_h \mathbf{b} / a}{\mathbf{b}' \mathbf{S}_e \mathbf{b} / \nu_e} \\ &= \frac{(\mathbf{b}' \widehat{\boldsymbol{\Theta}}) [\mathbf{C} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}']^{-1} (\widehat{\boldsymbol{\Theta} \mathbf{b}}) / a}{\mathbf{b}' \mathbf{U}' \mathbf{Y}' [\mathbf{I}_N - \mathbf{X} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'] \mathbf{Y} \mathbf{U} \mathbf{b} / \nu_e}. \end{aligned} \tag{16.30}$$

The distributions of \mathbf{S}_h and \mathbf{S}_e are both Wishart. Hence $q_j = \mathbf{b}' \mathbf{S}_j \mathbf{b}$ is distributed as a scaled chi square. The univariate UI test is

$$\begin{aligned} \widehat{H}_A(\mathbf{b}) &\equiv \widehat{H}_A(\mathbf{b}; \mathbf{Y}) \\ &= \mathbb{B}[F(\mathbf{b}; \mathbf{Y}) > f_{\text{crit}}], \end{aligned} \tag{16.31}$$

with decision function

$$\begin{aligned} \widehat{H}_0(\mathbf{b}) &\equiv \widehat{H}_0(\mathbf{b}; \mathbf{Y}) \\ &= \mathbb{B}[F(\mathbf{b}; \mathbf{Y}) \leq f_{\text{crit}}] \\ &= 1 - \widehat{H}_A(\mathbf{b}). \end{aligned} \tag{16.32}$$

The appropriate critical value, f_{crit} , is the $100(1 - \alpha)$ percentile of the $F(a, \nu_e)$ distribution. The appropriate critical value does *not* depend on \mathbf{b} .

Corresponding acceptance and rejection regions are

$$AR(H_0) = \bigcap_{\mathbf{b} \neq \mathbf{0}} AR[H_0(\mathbf{b})] \quad (16.33)$$

$$RR(H_0) = \bigcup_{\mathbf{b} \neq \mathbf{0}} RR[H_0(\mathbf{b})]. \quad (16.34)$$

The regions are sets in the sample space, $AR(H_0) = \{\mathbf{Y} : \hat{H}_A = \text{TRUE}\}$ and $RR(H_0) = \{\mathbf{Y} : \hat{H}_0 = \text{FALSE}\}$, while

$$\begin{aligned} AR[H_0(\mathbf{b})] &= \{\mathbf{Y} : \hat{H}_0(\mathbf{b}) = \text{TRUE}\} \\ &= \{\mathbf{Y} : F(\mathbf{b}; \mathbf{Y}) \leq f_{\text{crit}}\} \end{aligned} \quad (16.35)$$

and

$$\begin{aligned} RR[H_0(\mathbf{b})] &= \{\mathbf{Y} : \hat{H}_0(\mathbf{b}) = \text{FALSE}\} \\ &= \{\mathbf{Y} : F(\mathbf{b}; \mathbf{Y}) > f_{\text{crit}}\}. \end{aligned} \quad (16.36)$$

The decision function is $\hat{H}_0 = \mathbb{B}[\mathbf{Y} \in AR(H_0)]$ and the test is $\hat{H}_A = \mathbb{B}[\mathbf{Y} \in RR(H_0)]$. Thus by the union-intersection principle we have

$$\hat{H}_A = \bigvee_{\mathbf{b} \neq \mathbf{0}} \hat{H}_A(\mathbf{b}) \quad (16.37)$$

$$H_0 = \bigwedge_{\mathbf{b} \neq \mathbf{0}} H_0(\mathbf{b}). \quad (16.38)$$

With the above notation the heart of the proof is simple:

$$\begin{aligned} \hat{H}_0 &= \mathbb{B}[\mathbf{Y} \in AR(H_0)] \\ &= \bigwedge_{\mathbf{b} \neq \mathbf{0}} \mathbb{B}\{\mathbf{Y} \in AR[H_0(\mathbf{b})]\} \\ &= \bigwedge_{\mathbf{b} \neq \mathbf{0}} \mathbb{B}[F(\mathbf{b}; \mathbf{Y}) \leq f_{\text{crit}}] \\ &= \mathbb{B}\left[\sup_{\mathbf{b} \neq \mathbf{0}}\{F(\mathbf{b}; \mathbf{Y})\} \leq f_{\text{crit}}\right] = \mathbb{B}(\lambda \leq f_{\text{crit}}). \end{aligned} \quad (16.39)$$

Similar equalities allow proving $\hat{H}_A = \mathbb{B}(\lambda > f_{\text{crit}})$. Thus the UI statistic is

$$\begin{aligned} \lambda &= \sup_{\mathbf{b} \neq \mathbf{0}}\{F(\mathbf{b}; \mathbf{Y})\} \\ &= \sup_{\mathbf{b} \neq \mathbf{0}} \left(\frac{\mathbf{b}' \mathbf{S}_h \mathbf{b}}{\mathbf{b}' \mathbf{S}_e \mathbf{b}} \right) \left(\frac{\nu_e}{a} \right). \end{aligned} \quad (16.40)$$

Given \mathbf{S}_e^{-1} exists with probability 1, so does the Cholesky decomposition $\mathbf{S}_e^{-1} = \mathbf{L}'\mathbf{L}$ and \mathbf{L}^{-1} . If $\mathbf{w} = (\mathbf{L}^{-1})'\mathbf{b}$, then $\mathbf{b} = \mathbf{w}'\mathbf{L}$ and

$$\begin{aligned} \widehat{\phi} &= \sup_{\mathbf{w} \neq \mathbf{0}} \left(\frac{\mathbf{w}' \mathbf{L} \mathbf{S}_h \mathbf{L}' \mathbf{w}}{\mathbf{w}' \mathbf{w}} \right) \left(\frac{\nu_e}{a} \right) \\ &= \left(\frac{\nu_e}{a} \right) \widehat{\phi}_{\max}(\mathbf{L} \mathbf{S}_h \mathbf{L}'), \end{aligned} \tag{16.41}$$

with $\widehat{\phi}_{\max}$ the largest root of characteristic polynomial

$$\begin{aligned} |\mathbf{L} \mathbf{S}_h \mathbf{L}' - \widehat{\phi} \mathbf{I}| &\equiv |\mathbf{S}_h - \widehat{\phi} \mathbf{S}_e| \\ &\equiv |\mathbf{S}_h \mathbf{S}_e^{-1} - \widehat{\phi} \mathbf{I}| = 0. \end{aligned} \tag{16.42}$$

Thus the U-I test is

$$\begin{aligned} \widehat{H}_a &= \mathbb{B}(\widehat{\phi} > f_{\text{crit}}) \\ &= \mathbb{B}\left[(\nu_e/a) \widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1}) > f_{\text{crit}} \right] \\ &= \mathbb{B}\left[\widehat{\phi}_{\max}(\mathbf{S}_h \mathbf{S}_e^{-1}) > f_{\text{crit}} a/\nu_e \right]. \quad \square \end{aligned} \tag{16.43}$$

In practice, the statistic $\widehat{\rho}_1^2 = \widehat{\phi}_{\max}/(1 + \widehat{\phi}_{\max})$ (which is a monotone function of $\widehat{\phi}_{\max}$) is preferred because $0 \leq \widehat{\phi}_{\max} < \infty$ while $0 \leq \widehat{\rho}_1^2 \leq 1$. Therefore $\widehat{\rho}_1^2$ may be interpreted simply as the fraction of variance controlled. Since $\widehat{\phi}_{\max} = \widehat{\rho}_1^2/(1 - \widehat{\rho}_1^2)$, the critical values bear the relationship $l_\alpha = x_\alpha/(1 - x_\alpha)$.

The distribution of $\widehat{\rho}_1^2$ has been tabulated and charted by Heck, with the symbol θ_s . Heck's charts are given in the appendix of Morrison's (1990) book. Timm's book gives both tables and charts of the distribution. Also, Pillai (1956) gave an algorithm to approximate the right tail probability. Code for Pillai's approximation is listed in an appendix of Harris (1975). Free software LINMOD also implements the algorithm. The Appendix (Section A.2) of the present book contains a description of LINMOD and where to retrieve it from the Web. A missing value is returned if the p value exceeds 0.10 because the algorithm is guaranteed to work for small p values. Under H_0 , the distribution of $\widehat{\rho}_1^2$ has parameters $s = \min(a, b)$, $m = (|a-b|-1)/2$, and $n = (\nu_e - b - 1)/2$. Hypothesis H_0 is rejected if $\widehat{\rho}_1^2 > x_\alpha$.

Theorem 16.10 The substitution test (ST) principle leads to the following test for the GLM $_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with Gaussian errors, estimable $\Theta = \mathbf{C} \mathbf{B} \mathbf{U}$ ($a \times b$), $\text{rank}(\mathbf{X}) = r \leq q$, $\text{rank}(\mathbf{C}) = a \leq q$, $\text{rank}(\mathbf{U}) = b \leq p$, and known constants \mathbf{C} , \mathbf{U} , Θ_0 . The ST of $H_0 = \mathbb{B}(\Theta = \Theta_0)$ versus $H_A = -H_0$ is $\widehat{H}_A = \mathbb{B}[\text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1}) > t_0]$ with test statistic $\text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1})$, the Hotelling-Lawley trace, and appropriate critical value t_0 . Hypothesis H_0 is rejected for improbably large values of the test statistic relative to the appropriate critical value.

Proof. Arnold (1981, p. 363–364) included a proof.

The substitution principle may be described as the conditioning principle. If Σ_* were known, then the UMP, invariant, size- α test would reject the hypothesis for improbably large values of $\text{tr}(\mathcal{S}_h \Sigma_*^{-1})$. Substituting $\widehat{\Sigma}_*$ for Σ_* removes the conditioning and suggests a test which may or may not have desirable properties.

Under the null, the eigenvalues of $\mathcal{S}_h \mathcal{S}_t^{-1}$ will be close to zero. Data that do not support the null yield large eigenvalues. Hence each test statistic discussed can be thought of as a means of summarizing the information in the $s = \min(a, b)$ nonzero eigenvalues. Roy's criterion uses the largest eigenvalue, the LRT uses a product of nonzero eigenvalues (a one-to-one function of a geometric mean), and the two trace criteria look at weighted sums of the eigenvalues. Each statistic has a reasonable interpretation in terms of the eigenvalues.

The one- and two-sample versions of T^2 generalize the one sample and independent groups t tests to allow simultaneous consideration of p response variables. Both correspond to special cases of the four general statistics, which provide the same p value and conclusion in any situation with $a = 1$.

16.6 WHICH MULTIVARIATE TEST IS BEST?

The existence of four tests leads to an obvious question: Which one is best? Not surprisingly, the answer varies with the definition of "best" and the particular hypothesis being tested. Here we emphasize accuracy of test size, robustness to violation of assumptions, and power. The multipart answer to the question derives from a few dimensions and one set of eigenvalues for data analysis and a corresponding set of eigenvalues for power analysis.

All four test statistics are functions of the data solely through the $s = \min(a, b)$ nonzero eigenvalues of the $b \times b$ matrix $\widehat{\Omega} = \widehat{\Delta} \widehat{\Sigma}_*^{-1} = \mathcal{S}_h \mathcal{S}_e^{-1} \nu_e$. Equivalently, the s generalized, squared canonical correlations, $\{\widehat{\rho}_k^2\}$, the eigenvalues of $\widehat{\Delta}(\widehat{\Delta} + \nu_e \widehat{\Sigma}_*)^{-1} = \mathcal{S}_h(\mathcal{S}_h + \mathcal{S}_e)^{-1}$, suffice. The only additional values needed to compute any of the tests are the constants $\{a, b, \nu_e\}$ with $\nu_e = N - r$. Under the null, $\{a, b, \nu_e\}$ completely determine the distribution, and all four are exact size α .

Although all four tests are size α , typically they give four distinct p values. Reporting only the smallest is statistically dishonest due to the bias introduced. More precisely, the reported value gives a test with inflated test size.

Specifying the distributions of the four test statistics under the alternative requires knowing only $\{a, b, \nu_e\}$ and the nonzero eigenvalues of the noncentrality matrix $\Omega = \Delta \Sigma_*^{-1}$, or equivalently $\{\rho_k^2\}$, the eigenvalues of $\Delta(\Delta + \nu_e \Sigma_*)^{-1}$. Although $\text{rank}(\widehat{\Omega}) = s$ with probability 1 and $1 \leq s \leq b$, $\text{rank}(\Omega) = s_*$ with $0 \leq s_* \leq s$. Here $s_* = 0$ iff $\Omega = \mathbf{0}$ iff $\Theta = \Theta_0$ iff the null hypothesis is true.

If $s = 1$, then the four multivariate test statistics (1) are one to one functions of each other, and (2) provide exactly the same p value and power and, (3) the unique test is optimal in many ways. The test is then exactly size α , invariant, provides the likelihood ratio test and the union-intersection test, and is uniformly most powerful among all size- α and invariant tests.

If $s > 1$, no uniformly most powerful test exists among the class of size- α and invariant tests for finite samples. Which test is most powerful depends on the particular set of eigenvalues of Ω or equivalently the particular set of $\{\rho_k^2\}$. Some performance differences among the tests are known. Olson (1974, 1976, 1979) provided extensive studies of power and robustness for the multivariate tests. He characterized alternative hypotheses as involving either concentrated noncentrality if $s_* = 1$ or diffuse noncentrality if $s_* > 1$. Concentrated noncentrality ensures RLR is the most powerful test of the four, while diffuse noncentrality ensures RLR is the least powerful test of the four. Olson concluded RLR was the least robust while PBT was the most robust. For the range of cases he considered, he strongly preferred PBT, with LRT in second place.

Example 16.1 The following unusual example occurred in an actual data analysis. For a data set with $s = 5$, $m = 0$, $N = 10$, the following p values were obtained: Wilks lambda $p = 0.1896$, Pillai's trace $p = 0.7799$, Hotelling-Lawley trace $p = 0.0030$, Roy's largest root $p = 0.0001$. Such examples are uncommon. Usually the various test multivariate statistics are in approximate agreement.

16.7 UNIVARIATE APPROACH TO REPEATED MEASURES: UNIREP

Muller and Barton (1989) detailed the origin and history of the UNIREP tests. Although they are not invariant to a full-rank transformation of the transformed response data, they do meet a weaker criterion, namely invariance to a full-rank orthonormal transformation. All linearly invariant tests, including the four "MULTIREP" tests, are also orthonormal invariant.

The sphericity property is met whenever Σ_* has all eigenvalues of the same size. With sphericity, $\epsilon = 1$, and the uncorrected test provides an exact size α with uniformly most power, among the class of similarly invariant tests. Without sphericity, the uncorrected test typically has inflated test size, the Geisser-Greenhouse and Huynh-Feldt tests are approximately size α , and the Box test is nearly always very conservative.

John (1972) discussed the likelihood ratio test of sphericity (Mauchly's test) and the locally (near the null) most powerful size- α test of sphericity. The likelihood ratio statistic is a one-to-one function of $|\widehat{\Sigma}_*|/\text{tr}(\widehat{\Sigma}_*)$, while the locally most powerful test is a one-to-one function of $\widehat{\epsilon} = b^{-1}\text{tr}(\widehat{\Sigma}_*)/\text{tr}(\widehat{\Sigma}_*^2)$. In our opinion, neither should be used as a "screening test" to choose a UNIREP test, for the same reasons a test of heterogeneity should not be used to choose between an unadjusted or Satterthwaite t test. The locally most powerful nature of the sphericity test based on $\widehat{\epsilon}$ leads to the speculation that the Geisser-Greenhouse test shares the property.

Coffey and Muller (2003) proved the following lemma, which characterizes the noncentral distribution of the UNIREP statistic. The lemma uses the spectral decomposition $\Sigma_* = \Upsilon \text{Dg}(\lambda) \Upsilon'$, with $\Upsilon' \Upsilon = I_b$, which allows defining $\Delta_\Upsilon = \Upsilon' \Delta \Upsilon$ and $\Omega_\Upsilon = \Delta_\Upsilon \text{Dg}(\lambda)^{-1}$.

Lemma 16.2 Except for known constants, the distribution of the UNIREP test statistic, $T_u = [\text{tr}(\widehat{\Delta})/a]/\text{tr}(\widehat{\Sigma}_*)$, is completely and exactly determined by the b eigenvalues of Σ_* , λ , and the b noncentralities $\omega_\Upsilon = \{\omega_{\Upsilon kk}\} = \{\mathbf{v}'_k(\Theta - \Theta_0)'M^{-1}(\Theta - \Theta_0)\mathbf{v}_k/\lambda_k\}$, the diagonal elements of $\Omega_\Upsilon = \Delta_\Upsilon \text{Dg}(\lambda)^{-1}$. The same characterization holds for all four UNIREP tests.

Proof. Coffey and Muller (2003) provided a proof.

It is straightforward to prove that the eigenvalues of Ω_Υ and Ω coincide. Of course, the eigenvalues of Ω_Υ coincide with the diagonal values only when Ω_Υ is diagonal. Hence it is not surprising that Monte Carlo simulations demonstrate the UNIREP tests may be more or less powerful than any of the MULTIREP tests, depending on λ , ω_Υ , and the eigenvalues of Ω_Υ , which are equivalent to $\{\rho_k^2\}$.

16.8 MORE ON INVARIANCE PROPERTIES

Theorem 16.11 (a) The eigenvalues of $S_h S_e^{-1}$, (b) the canonical correlations, (c) the usual four MULTIREP test statistics, and (d) the associated p values are invariant to a square and full-rank transformation simultaneously applied to $Y_U = YU$ and Θ_0 , with $Y_T = Y_U T$ and $\Theta_{0T} = \Theta_0 T$. The same invariance holds under the same type of transformation to U and Θ_0 .

Proof. Here $S_e = FF'$ and $S_e^{-1} = F^{-t}F^{-1}$. A similarity transformation allows proving the eigenvalues of $S_h S_e^{-1} = S_h F^{-t}F^{-1}$ coincide with the eigenvalues of $F^{-1}S_h F^{-t}$. If $H = X(X'X)^{-1}X'$ and $M = C(X'X)C'$, then

$$\begin{aligned} S_e &= U'\widehat{\Sigma}U\nu_e \\ &= U'[Y'(I - H)Y/\nu_e]U\nu_e \\ &= U'Y'(I - H)YU = Y'_U A_e Y_U \end{aligned} \quad (16.44)$$

and

$$\begin{aligned} S_h &= (\widehat{\Theta} - \Theta_0)'M^{-1}(\widehat{\Theta} - \Theta_0) \\ &= (C\widehat{B}U - \Theta_0)'M^{-1}(C\widehat{B}U - \Theta_0) \\ &= [C(X'X)^{-1}X'YU - \Theta_0]'M^{-1}[C(X'X)^{-1}X'YU - \Theta_0] \\ &= [Y'_U X(X'X)^{-1}C' - \Theta'_0]M^{-1}[C(X'X)^{-1}X'Y_U - \Theta_0]. \end{aligned} \quad (16.45)$$

If $\Theta_0 = \mathbf{0}$, then

$$\begin{aligned} S_h &= Y'_U X(X'X)^{-1}C'M^{-1}C(X'X)^{-1}X'Y_U \\ &= Y'_U A_h Y_U. \end{aligned} \quad (16.46)$$

In turn,

$$\begin{aligned}
 (\mathbf{S}_h \mathbf{S}_e^{-1} - \lambda \mathbf{I}) \mathbf{x} &= \mathbf{0} \\
 [(\mathbf{Y}'_U \mathbf{A}_h \mathbf{Y}_U)(\mathbf{Y}'_U \mathbf{A}_e \mathbf{Y}_U)^{-1} - \lambda \mathbf{I}] \mathbf{x} &= \mathbf{0} \\
 [(\mathbf{Y}'_U \mathbf{A}_h \mathbf{Y}_U) \mathbf{T} \mathbf{T}^{-1} (\mathbf{Y}'_U \mathbf{A}_e \mathbf{Y}_U)^{-1} \mathbf{T}^{-t} \mathbf{T}' - \lambda \mathbf{T}^{-t} \mathbf{T}'] \mathbf{x} &= \mathbf{0} \\
 [(\mathbf{Y}'_U \mathbf{A}_h \mathbf{Y}_U \mathbf{T})(\mathbf{T}' \mathbf{Y}'_U \mathbf{A}_e \mathbf{Y}_U \mathbf{T})^{-1} - \lambda \mathbf{T}^{-t}] \mathbf{T}' \mathbf{x} &= \mathbf{0} \\
 \mathbf{T}' [(\mathbf{Y}'_U \mathbf{A}_h \mathbf{Y}_U \mathbf{T})(\mathbf{T}' \mathbf{Y}'_U \mathbf{A}_e \mathbf{Y}_U \mathbf{T})^{-1} - \lambda \mathbf{T}^{-t}] \mathbf{T}' \mathbf{x} &= \mathbf{T}' \mathbf{0} \\
 [(\mathbf{T}' \mathbf{Y}'_U \mathbf{A}_h \mathbf{Y}_U \mathbf{T})(\mathbf{T}' \mathbf{Y}'_U \mathbf{A}_e \mathbf{Y}_U \mathbf{T})^{-1} - \lambda \mathbf{I}] \mathbf{x}_T &= \mathbf{0} \\
 [(\mathbf{Y}'_T \mathbf{A}_h \mathbf{Y}_T)(\mathbf{Y}'_T \mathbf{A}_e \mathbf{Y}_T)^{-1} - \lambda \mathbf{I}] \mathbf{x}_T &= \mathbf{0}. \quad (16.47)
 \end{aligned}$$

The last form demonstrates the eigenvalues coincide (although eigenvectors differ by the factor of \mathbf{T}') if $\Theta_0 = \mathbf{0}$. Allowing $\Theta_0 \neq \mathbf{0}$ is covered below. \square

Corollary 16.11 Without loss of generality, analysis or simulation of the MULTIREP tests for any $\{\Sigma_{*1} = \Phi \Phi', \Delta_1, \Omega_1 = \Delta_1 \Sigma_{*1}^{-1}\}$ can be based on the equivalent model with $\Sigma_{*2} = \mathbf{I}_b$, $\Delta_2 = \Phi^{-1} \Delta_1 \Phi^{-t}$, and $\Omega_2 = \Delta_2 \Sigma_{*2}^{-1} = \Delta_2$. Also, $\Omega_2 = \Omega'_2$ has the same eigenvalues as Ω_1 (which need not be symmetric).

Proof. If $\Sigma_{*1} = \Phi \Phi'$ and $\mathbf{T} = \Phi^{-t}$, then $[\text{row}_i(\mathbf{Y}_T)]' \sim \mathcal{N}_b[(\mathbf{X}_i \mathbf{B} \mathbf{U} \Phi^{-t})', \mathbf{I}_b]$. Also, $\hat{\Delta}_1 \sim \mathcal{W}_b(a, \Sigma_*, \Delta_1)$ and $\Phi^{-1} \hat{\Delta}_1 \Phi^{-t} \sim \mathcal{W}_b(a, \mathbf{I}_b, \Phi^{-1} \Delta_1 \Phi^{-t})$. Finally, Ω_1 is similar to $\Phi^{-1}(\Omega_1)\Phi = \Phi^{-1}(\Delta_1 \Phi^{-t} \Phi^{-1})\Phi = \Phi^{-1} \Delta_1 \Phi^{-t} = \Delta_2 = \Omega_2$ and therefore has the same eigenvalues. \square

Theorem 16.12 (a) The matrices \mathbf{S}_e and \mathbf{S}_h are invariant to a square, full-rank, simultaneous transformation of the rows of \mathbf{C} and Θ_0 , with $\mathbf{C}_T = \mathbf{T} \mathbf{C}$ and $\Theta_{0T} = \mathbf{T} \Theta_0$.

(b) The eigenvalues of $\mathbf{S}_h \mathbf{S}_e^{-1}$, the canonical correlations, the usual four MULTIREP test statistics, and their associated p values are all invariant in the same way.

Proof. **(a)** \mathbf{S}_e is not a function of \mathbf{C} . If \mathbf{T} is $a \times a$ and full rank, with $\mathbf{C}_T = \mathbf{T} \mathbf{C}$,

$$\begin{aligned}
 \mathbf{S}_h &= (\mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \Theta_0)' [\mathbf{C} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}']^{-1} (\mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \Theta_0) \quad (16.48) \\
 &= (\mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \Theta_0)' (\mathbf{T}^{-1} \mathbf{T})' [\mathbf{C} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}']^{-1} \mathbf{T}^{-1} \mathbf{T} (\mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \Theta_0) \\
 &= (\mathbf{T}' \mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \mathbf{T} \Theta_0)' [\mathbf{T} \mathbf{C} (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}' \mathbf{T}']^{-1} (\mathbf{T} \mathbf{C} \tilde{\mathbf{B}} \mathbf{U} - \mathbf{T} \Theta_0) \\
 &= (\mathbf{C}_T \tilde{\mathbf{B}} \mathbf{U} - \Theta_{0T})' [\mathbf{C}_T (\mathbf{X}' \mathbf{X})^{-1} \mathbf{C}'_T]^{-1} (\mathbf{C}_T \tilde{\mathbf{B}} \mathbf{U} - \Theta_{0T}),
 \end{aligned}$$

which suffices to prove the invariance of \mathbf{S}_h . Part **(b)** follows immediately. \square

As noted earlier, the parameters determining the distributions of the secondary parameter estimators make it obvious that any reasonable test statistic must be a function of $\hat{\Delta}$ and $\hat{\Sigma}_*$. Given that the eigenvalues of $\mathbf{S}_h \mathbf{S}_e^{-1}$ are invariant, how do

we know they are the only invariants for a testable hypothesis? We have proven all four statistics depend on the data *only* through the eigenvalues. Therefore any other feature may change and not affect the test, unless the change carries over into the eigenvalues. They are also obviously the minimal sufficient statistics because changing any one of them changes the test statistics (except perhaps for RLR).

Theorem 16.13 If the model spans an intercept ($\mathbf{X}\mathbf{t}_0 = \mathbf{1}_N$) and the hypothesis excludes the intercept ($\mathbf{C}\mathbf{t}_0 = \mathbf{0}$), then the following hold.

- (a) Matrices \mathbf{S}_h and \mathbf{S}_e , the eigenvalues of $\mathbf{S}_h\mathbf{S}_e^{-1}$, and the canonical correlations, are invariant to a location shift of the form $\mathbf{Y}_U + \mathbf{1}_N\mathbf{t}'_1$ for $b \times 1$ constant vector \mathbf{t}_1 and $\mathbf{Y}_U = \mathbf{Y}\mathbf{U}$.
- (b) The MULTIREP test statistics and associated p values are similarly invariant.
- (c) The UNIREP test statistics and associated p values are similarly invariant.

Proof. With $\mathbf{H} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$, the matrices \mathbf{S}_h and \mathbf{S}_e and the test statistics depend on \mathbf{Y}_U only through

$$\hat{\boldsymbol{\Theta}} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}_U \quad (16.49)$$

$$\hat{\mathbf{E}}_U = (\mathbf{I}_N - \mathbf{H})\mathbf{Y}_U, \quad (16.50)$$

with $\hat{\boldsymbol{\Sigma}}_* = \hat{\mathbf{E}}_U'\hat{\mathbf{E}}_U/\nu_e$. Writing

$$\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'(\mathbf{Y}_U + \mathbf{1}_N\mathbf{t}'_1) = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y}_U + \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{1}_N\mathbf{t}'_1 \quad (16.51)$$

$$(\mathbf{I}_N - \mathbf{H})(\mathbf{Y}_U + \mathbf{1}_N\mathbf{t}'_1) = (\mathbf{I}_N - \mathbf{H})\mathbf{Y}_U + (\mathbf{I}_N - \mathbf{H})\mathbf{1}_N\mathbf{t}'_1 \quad (16.52)$$

allows concluding location invariance reduces to proving

$$\mathbf{0} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{1}_N\mathbf{t}'_1 \quad (16.53)$$

$$\mathbf{0} = [\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{1}_N\mathbf{t}'_1. \quad (16.54)$$

The zero matrices of the two equations have dimensions $a \times b$ and $N \times b$. In turn, it suffices to prove

$$\mathbf{0} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{1}_N \quad (16.55)$$

$$\mathbf{0} = [\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{1}_N. \quad (16.56)$$

The zero matrices have dimensions $a \times 1$ and $n \times 1$.

In the simplest case, \mathbf{X} is full rank, the first column is $\mathbf{1}_N$, and the remaining columns are centered. Therefore

$$\mathbf{X}'\mathbf{1}_N = \begin{bmatrix} N \\ \mathbf{0}_{(q-1) \times 1} \end{bmatrix} \quad (16.57)$$

and

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} N & \mathbf{0} \\ \mathbf{0} & \mathbf{X}'_d\mathbf{X}_d \end{bmatrix}^{-1} = \begin{bmatrix} N^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}'_d\mathbf{X}_d)^{-1} \end{bmatrix}. \quad (16.58)$$

If the test excludes the intercept, then $\mathbf{C} = [\mathbf{0}_{a \times 1} \quad \mathbf{C}_d]$ and

$$\begin{aligned} C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{1}_N &= [\mathbf{0}_{a \times 1} \quad C_d] \begin{bmatrix} N^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}'_d\mathbf{X}_d)^{-1} \end{bmatrix} \begin{bmatrix} N \\ \mathbf{0}_{(q-1) \times 1} \end{bmatrix} \\ &= [\mathbf{0}_{a \times 1} \quad C_d] \begin{bmatrix} 1 \\ \mathbf{0}_{(q-1) \times 1} \end{bmatrix} = \mathbf{0}_{a \times 1}. \end{aligned} \tag{16.59}$$

Furthermore

$$\begin{aligned} [\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{1}_N &= \mathbf{1}_N - \mathbf{X} \begin{bmatrix} N^{-1} & \mathbf{0} \\ \mathbf{0} & (\mathbf{X}'_d\mathbf{X}_d)^{-1} \end{bmatrix} \begin{bmatrix} N \\ \mathbf{0}_{(q-1) \times 1} \end{bmatrix} \\ &= \mathbf{1}_N - \mathbf{X} \begin{bmatrix} 1 \\ \mathbf{0}_{(q-1) \times 1} \end{bmatrix} \\ &= \mathbf{1}_N - \mathbf{1}_N = \mathbf{0}_{N \times 1}, \end{aligned} \tag{16.60}$$

which completes the proof of the special case. The theory of linearly equivalent models allows extending the result to estimable parameters in LTFR cases. \square

Theorem 16.14 The UNIREP test statistic, for any population covariance pattern, is invariant to a FR transformation of columns of \mathbf{U} and Θ_0 by orthonormal \mathbf{T} (but not general FR \mathbf{T}). Here $\mathbf{T}^{-1} = \mathbf{T}'$ and hence $\mathbf{T}\mathbf{T}' = \mathbf{I}_b$.

Proof.

$$F_u = \frac{\text{tr}(\mathbf{S}_h)/(ab)}{\text{tr}(\mathbf{S}_e)/(b\nu_e)} = \frac{\text{tr}(\mathbf{S}_h\mathbf{T}\mathbf{T}')/(ab)}{\text{tr}(\mathbf{S}_e\mathbf{T}\mathbf{T}')/(b\nu_e)} = \frac{\text{tr}(\mathbf{T}'\mathbf{S}_h\mathbf{T})/(ab)}{\text{tr}(\mathbf{T}'\mathbf{S}_e\mathbf{T})/(b\nu_e)}. \quad \square$$

Corollary 16.14 In a simulation, choosing $\mathbf{T} = \mathbf{Y}$, with $\Sigma_* = \mathbf{Y}\text{Dg}(\lambda)\mathbf{Y}'$ and $\mathbf{Y}_T = \mathbf{Y}\mathbf{U}\mathbf{T}$, gives $[\text{row}_i(\mathbf{Y}_T)]' \sim \mathcal{N}_b[(\mathbf{X}_i\mathbf{B}\mathbf{U}\mathbf{Y})', \text{Dg}(\lambda)]$.

Theorem 16.15 Without loss of generality, $\Theta_0 = \mathbf{0}$ may be assumed.

Proof. A multivariate GLM $_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors, fixed \mathbf{X} , $N \gg r = \text{rank}(\mathbf{X})$, has model equation $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$, with testable secondary parameter $\Theta = \mathbf{C}\mathbf{B}$ and corresponding hypothesis $H_0 : \Theta = \Theta_0$. Choosing

$$\mathbf{Y}_T = \mathbf{Y} - \mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\Theta_0 \tag{16.61}$$

defines the model

$$\mathbf{Y}_T = \mathbf{X}\mathbf{B}_T + \mathbf{E}_T, \tag{16.62}$$

with associated estimable secondary parameters $\Theta_T = \mathbf{C}\mathbf{B}_T$ and corresponding

hypothesis $H_0 : \Theta_T = \mathbf{0}$. Estimability ensures $C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = C$. Therefore

$$\begin{aligned}\hat{\Theta}_T &= C\tilde{B}_T \\ &= C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'[\mathbf{Y} - \mathbf{X}C'(CC')^{-1}\Theta_0] \\ &= C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} - [C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}]C'(CC')^{-1}\Theta_0 \\ &= C(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{Y} - [C]C'(CC')^{-1}\Theta_0 \\ &= \hat{\Theta} - \Theta_0.\end{aligned}\tag{16.63}$$

□

The result holds for hypothesis testing, no matter what kind of predictors. It is unclear whether the result applies to power analysis for random predictors.

Theorem 16.16 For any testable general linear hypothesis for a GLM, a linearly equivalent model may be found in which $C = [I_a \ \mathbf{0}]$ provides the original test.

Proof. A GLM (with any sort of predictors), with $N \gg r = \text{rank}(\mathbf{X})$, has model equation $\mathbf{Y} = \mathbf{X}\mathbf{B} + \mathbf{E}$, with testable secondary parameter $\Theta = C\mathbf{B}$ and corresponding hypothesis $H_0 : \Theta = \mathbf{0}$. For any C of dimension $a \times q$, rank a , the singular value decomposition allows writing $C = L[\text{Dg}(s_C) \ \mathbf{0}_{a \times (q-a)}]R'$ with L and $\text{Dg}(s_C)$ $a \times a$ and R $q \times q$. In turn, $CC' = L\text{Dg}(s_C^2)L'$, and $LL' = L'L = I_a$. With R_1 $q \times a$ and R_0 $q \times (q - a)$,

$$C'C = \begin{bmatrix} R_1 & R_0 \end{bmatrix} \begin{bmatrix} \text{Dg}(s_C^2) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} R'_1 \\ R'_0 \end{bmatrix},\tag{16.64}$$

with $RR' = R'R = I_q$. If $C_\perp = R'_0$, which is $(q - a) \times q$, then

$$\begin{aligned}T &= \begin{bmatrix} C \\ R'_0 \end{bmatrix} \\ &= \begin{bmatrix} L\text{Dg}(s_C)R'_1 \\ R'_0 \end{bmatrix} \\ &= \begin{bmatrix} L & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \begin{bmatrix} \text{Dg}(s_C) & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} R'\end{aligned}\tag{16.65}$$

and

$$\begin{aligned}T^{-1} &= \begin{bmatrix} R_1 & R_0 \end{bmatrix} \begin{bmatrix} \text{Dg}(s_C)^{-1} & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \begin{bmatrix} L' & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \\ &= \begin{bmatrix} R_1 & R_0 \end{bmatrix} \begin{bmatrix} \text{Dg}(s_C)^{-1}L' & \mathbf{0} \\ \mathbf{0} & I \end{bmatrix} \\ &= \begin{bmatrix} R_1\text{Dg}(s_C)^{-1}L' & R_0 \end{bmatrix} \\ &= \begin{bmatrix} C^+ & R_0 \end{bmatrix} \\ &= \begin{bmatrix} C'(CC')^{-1} & R_0 \end{bmatrix}.\end{aligned}\tag{16.66}$$

In turn,

$$\begin{aligned}
 \mathbf{Y} &= \mathbf{X}\mathbf{T}^{-1}\mathbf{T}\mathbf{B} + \mathbf{E} \\
 &= [\mathbf{X}\mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1} \quad \mathbf{X}\mathbf{R}\mathbf{R}_0] \begin{bmatrix} \boldsymbol{\Theta} \\ \boldsymbol{\Theta}_\perp \end{bmatrix} + \mathbf{E} \\
 &= [\mathbf{X}_{Tc} \quad \mathbf{X}_{T0}] \begin{bmatrix} \boldsymbol{\Theta} \\ \boldsymbol{\Theta}_\perp \end{bmatrix} + \mathbf{E} \\
 &= \mathbf{X}_T \begin{bmatrix} \boldsymbol{\Theta} \\ \boldsymbol{\Theta}_\perp \end{bmatrix} + \mathbf{E}.
 \end{aligned}
 \tag{16.67}$$

□

The proof is a special case of one approach to a “completion” theorem, seen in various contexts in linear models and matrix theory. The basic idea is part of the theory and process behind the concept of linearly equivalent models. A parallel approach is useful for identifying the estimable part of a LTFR model.

A similar result may be proven for $\mathbf{U} = \mathbf{I}$. The result holds for hypothesis testing, no matter whether predictors are fixed or only conditionally fixed. It is unclear whether the result applies to random predictors.

16.9 TESTS OF HYPOTHESES ABOUT Σ

Multivariate models raise the interesting possibility of estimating constrained structures and testing hypotheses about the structure of Σ . Many exact and approximate results have been developed for such tasks. Such techniques may be roughly described as falling into one of four types.

One class of tests concerns testing hypotheses about functions of Σ , such as the trace or determinant. For central Wishart matrices, both have simple forms. Each can be thought of as the generalized variance, with the trace corresponding to the arithmetic average eigenvalue and the determinant to the geometric mean,

$$\text{tr}(\Sigma)/p = \sum_{k=1}^p \lambda_k/p = \bar{\lambda}_A \tag{16.68}$$

$$|\Sigma|^{1/p} = \left(\prod_{k=1}^p \lambda_k \right)^{1/p} = \bar{\lambda}_G. \tag{16.69}$$

The interpretation is strengthened by recognizing that the eigenvalues of a covariance matrix equal the variances of the underlying principal components.

A second class of tests involves $H_0 : \Sigma = \Sigma_0$, based on $\widehat{\Sigma} \sim \mathcal{W}_p(\nu, \Sigma_0, \mathbf{0})$ and Σ_0 a particular structure. Anderson (2004), Morrison (1990), and Timm (2002, Section 3.8) provided useful treatments of some techniques. Some specific tests have been developed for special patterns, including $\Sigma_0 = \sigma^2 \mathbf{I}_p$, sphericity, and $\Sigma_0 = \sigma^2[\rho \mathbf{1}_p \mathbf{1}'_p + (1 - \rho) \mathbf{I}_p]$, compound symmetry (Lemma 1.33).

Tests for sphericity have an unfortunate history in the “univariate approach” to repeated measures. Violation of sphericity may greatly inflate the test size of the uncorrected UNIREP test. The problem led to using a nonsignificant test of sphericity as justification for using an uncorrected test. The process is the multivariate generalization of testing for homogeneity of variance between groups as part of a univariate ANOVA for independent observations. The process fails to control test size while providing a false sense of security.

In contrast, the following decision path is valid. One should assume compound symmetry if and only if the sampling scheme guarantees it. Most importantly, we have never encountered repeated measures in time which seemed likely to meet the assumption. If the sampling scheme does imply meeting the assumption, then use the uncorrected test. Otherwise, use either the Geisser-Greenhouse corrected test (preferred by the present authors) or the Huynh-Felt corrected test, which is likely somewhat more widely used. Muller, Edwards, Simpson, and Taylor (2006) and Muller and Barton (1989) provided additional discussion. Careful consideration of Kirk's (1995) recommendation for a three-step process allows concluding it always gives the same test result as always using the Geisser-Greenhouse test.

A third class of tests involve techniques known as *factor analysis*, which provide tools for developing models of correlation and covariance matrices. The key idea is to assume each response variable equals a linear combination of a set of unobserved, underlying, *latent variables* plus a component unique to the variable. The name *common factor model* reflects the decomposition in terms of shared, or common, and unique pieces. Advances in computing power allowed focusing on maximum likelihood methods for *covariance structure* models. From the perspective of the multivariate GLM, such a model allows only a simple design matrix, such as one between group factor, and complex specification of covariance structure. McDonald (1985) and Jöreskog (1993) gave thorough presentations.

The fourth and final type of tests involves mixed models. A general linear mixed model requires an explicit choice for a covariance model, commonly described as specifying the random effects. The validity and numerical feasibility of such a model depends on choosing a scientifically defensible covariance model with sufficiently few parameters. Current methods use large-sample theory for comparing covariance structures. The fragility of mixed models with respect to the accuracy of the covariance model makes such statistics, especially in small samples, important in fitting mixed models. Regrettably, even less is known about statistics focused on such random effects than is known about statistics focused on fixed effects.

16.10 CONFIDENCE REGIONS FOR Θ

Section 2.10 contains a definition of confidence regions for θ in the univariate model. Inverting a hypothesis test creates a confidence region, and inverting a confidence region yields a unique hypothesis test. Confidence regions exist only

for parameters that are estimable and testable. Section 15.6 contains proofs of the results as well as a description of the inversion process and examples.

In Section 3.10, replacing $\{\mathbf{y}, \boldsymbol{\theta}, \sigma^2\}$ by $\{\text{vec}(\mathbf{Y}), \text{vec}(\boldsymbol{\Theta}), \text{vec}(\boldsymbol{\Sigma})\}$ allowed extending the definition of confidence regions from the univariate model to the multivariate model. Consequently, separate multivariate proofs are not needed for the following two theorems. In multivariate notation, the essence of the relationship between hypothesis tests and confidence regions is that $R(\mathbf{Y}) = [\boldsymbol{\Theta}_0 : \mathbf{Y} \in AR(\boldsymbol{\Theta}_0)]$ while $AR(\boldsymbol{\Theta}_0) = [\mathbf{Y}_* : \boldsymbol{\Theta}_0 \in R(\mathbf{Y}_*)]$.

Theorem 16.17 If for any $\boldsymbol{\Theta}_0 \in S$ a size- α test exists, $\phi(\mathbf{Y})$, of hypothesis $H(\boldsymbol{\Theta}_0) = \mathbb{B}(\boldsymbol{\Theta} = \boldsymbol{\Theta}_0)$, then there exists a corresponding confidence region for $\boldsymbol{\Theta}$ with confidence coefficient $c(\alpha) = 1 - \alpha$. Furthermore, if the acceptance region of $\phi(\mathbf{Y})$ is $AR(\boldsymbol{\Theta}_0) = [\mathbf{Y}_* : \phi(\mathbf{Y}_*) = 0]$, then the corresponding confidence region is $R(\mathbf{Y}) = [\boldsymbol{\Theta}_0 : \mathbf{Y} \in AR(\boldsymbol{\Theta}_0)]$.

Theorem 16.18 If there exists an exact $100(1 - \alpha)$ percent confidence region for $\boldsymbol{\Theta}$, then (a) a corresponding test procedure of size α exists for testing $H(\boldsymbol{\Theta}_0) = \mathbb{B}(\boldsymbol{\Theta} = \boldsymbol{\Theta}_0)$ for any $\boldsymbol{\Theta}_0 \in \boldsymbol{\Theta}$ and (b) $H(\boldsymbol{\Theta}_0)$ is a testable hypothesis. (c) If the confidence region is $R(\mathbf{Y})$, then the acceptance region of the corresponding test is $AR(\boldsymbol{\Theta}_0) = [\mathbf{Y}_* : \boldsymbol{\Theta}_0 \in R(\mathbf{Y}_*)]$.

Example 16.2 A $GLM_{N,p,q}(\mathbf{Y}; \mathbf{X}, \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors has $\text{rank}(\mathbf{X}) = r \leq q$, while $t(\mathbf{Y}; \boldsymbol{\Theta}_0)$ represents any of the multivariate test statistics based on $\mathbf{S}_h = (\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)' \mathbf{M}^{-1} (\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)$ and $\mathbf{S}_e = \nu_e \mathbf{U}' \hat{\boldsymbol{\Sigma}} \mathbf{U}$. When inverting the test based on t to obtain a confidence region for $\boldsymbol{\Theta}$ ($a \times b$), the test $\phi(\mathbf{Y}) = \mathbb{B}(t > t_{\text{crit}})$ of $H(\boldsymbol{\Theta}_0) = \mathbb{B}(\boldsymbol{\Theta} = \boldsymbol{\Theta}_0)$ has acceptance region

$$AR(\boldsymbol{\Theta}_0) = [\mathbf{Y} : t(\mathbf{Y}; \boldsymbol{\Theta}_0) \leq t_{\text{crit}}]. \tag{16.70}$$

Here t_{crit} is the appropriate critical value. The corresponding $100(1 - \alpha)$ percent ellipsoidal confidence region is

$$R(\mathbf{Y}) = [\boldsymbol{\Theta}_0 : \mathbf{Y} \in AR(\boldsymbol{\Theta}_0)]. \tag{16.71}$$

The region is a hyperellipsoidal and has ab dimensions. If $a = 1$, then

$$\begin{aligned} t(\mathbf{Y}; \boldsymbol{\Theta}_0) &= \text{tr}(\mathbf{S}_h \mathbf{S}_e^{-1}) \\ &= \text{tr}\{[\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}']^{-1}\} (\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0) \mathbf{S}_e^{-1} (\hat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)' \\ &= F(\mathbf{Y}; \boldsymbol{\Theta}_0) \frac{p(N - r)}{N(N - p)}, \end{aligned} \tag{16.72}$$

with $F(\mathbf{Y}; \boldsymbol{\Theta}_0) \sim F(p, n - p)$. It follows that the $100(1 - \alpha)$ percent ellipsoidal confidence region is

$$R(\mathbf{Y}) = [\boldsymbol{\Theta}_0 : (\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)\mathbf{S}_e^{-1}(\widehat{\boldsymbol{\Theta}} - \boldsymbol{\Theta}_0)' \leq \lambda], \quad (16.73)$$

in which $\lambda = t_{\text{crit}}/\text{tr}\{\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'\}^{-1}$.

EXERCISES

16.1 Model 1 is

$$\begin{aligned} [\mathbf{y}_1 \quad \mathbf{y}_2 \quad \mathbf{y}_3] &= \left(\mathbf{1}_n \otimes \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} \right) \mathbf{B}_1 + \mathbf{E} \\ \mathbf{Y} &= \mathbf{X}_1 \mathbf{B}_1 + \mathbf{E}. \end{aligned}$$

Finding a valid test of the general linear hypothesis $H_0 : \mathbf{C}\mathbf{B}\mathbf{U} = \boldsymbol{\Theta}_0$ requires verifying that (a) $\boldsymbol{\Theta} = \mathbf{C}\mathbf{B}\mathbf{U}$ is estimable and (b) the hypothesis is testable. In the following, with $\boldsymbol{\Theta}_j = \mathbf{C}_j \mathbf{B} \mathbf{U}_j$, verify (a) holds or does not hold, and (b) holds or does not hold. It is not sufficient to merely state the correct answer. You must briefly justify each positive or negative answer. For each estimable $\boldsymbol{\Theta}_j$, describe each element very briefly as a function of cell means.

$$16.1.1 \mathbf{C}_1 = [1 \quad 0 \quad 0 \quad 0]$$

$$\mathbf{U}_1 = \mathbf{I}_3$$

$$16.1.2 \mathbf{C}_2 = [0 \quad 1 \quad 0 \quad 0]$$

$$\mathbf{U}_2 = [\mathbf{1}_2 \quad -\mathbf{I}_2]'$$

$$16.1.3 \mathbf{C}_3 = [1 \quad 1 \quad 0 \quad 0]$$

$$\mathbf{U}_3 = [\mathbf{U}_2 \quad -\mathbf{U}_2]$$

$$16.1.4 \mathbf{C}_4 = [3 \quad 1 \quad 1 \quad 1]/3$$

$$\mathbf{U}_4 = \mathbf{I}_3$$

$$16.1.5 \mathbf{C}_5 = \begin{bmatrix} 0 & 1 & -1 & 0 \\ 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

$$\mathbf{U}_5 = \begin{bmatrix} 1 & -1 & 1 \\ 1 & 0 & -2 \\ 1 & 1 & 1 \end{bmatrix}$$

$$16.1.6 \mathbf{C}_6 = \begin{bmatrix} 0 & 1 & 0 & -1 \\ 0 & 0 & 1 & -1 \end{bmatrix}$$

$$\mathbf{U}_6 = \mathbf{U}_5$$

16.2 For this exercise, you may use any results in the book up to immediately before Lemma 16.1.

16.2.1 Prove part (a) of Lemma 16.1.

16.2.2 Prove part (b) of Lemma 16.1.

16.2.3 Prove part (c) of Lemma 16.1.

16.3.1 With \mathbf{S}_h as defined in Chapter 3 and Chapter 16, prove directly that $\text{tr}(\mathbf{S}_h)$ is invariant to a full-rank and orthonormal transformation of the columns of \mathbf{U} . Assume that $\boldsymbol{\Theta}_0 = \mathbf{0}$. (A separate proof, which is *not* part of this exercise, allows concluding that this may be done without loss of generality.)

16.3.2 With $\widehat{\boldsymbol{\Sigma}}_*$ as defined in Chapter 3 and Chapter 16, prove directly that $\text{tr}(\widehat{\boldsymbol{\Sigma}}_*)$ is invariant to a full rank and orthonormal transformation of the columns of \mathbf{U} .

16.3.3 Prove directly that the UNIREP F statistic is invariant to multiplying \mathbf{C} (between-subject contrast matrix) by a nonzero constant. Assume that $\boldsymbol{\Theta}_0 = \mathbf{0}$.

16.3.4 Prove directly that the UNIREP F statistic is invariant to multiplying \mathbf{U} (within-subject contrast matrix) by a nonzero constant. Assume that $\boldsymbol{\Theta}_0 = \mathbf{0}$.

16.4 A complete, unbalanced, repeated-measures ANOVA, with p repeated measures, may be written using cell mean coding as

$$Y = \begin{bmatrix} \mathbf{1}_{N_1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \ddots & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{1}_{N_G} \end{bmatrix} \begin{bmatrix} \mu_{1,1} & \cdots & \mu_{1,p} \\ \vdots & \ddots & \vdots \\ \mu_{G,1} & \cdots & \mu_{G,p} \end{bmatrix} + E.$$

Assuming i.i.d. $[\text{row}_i(E)]' \sim \mathcal{N}_p(\mathbf{0}, \Sigma)$, a testable general linear hypothesis has the form $H_0 : \Theta = \Theta_0$, with $\Theta = CBU$. It is convenient to write $\sum_{g=1}^G N_g = N_+$.

16.4.1 Express the distribution of Y in terms of a (direct-product) matrix Gaussian.

16.4.2 Ignoring group, the set of p averages (grand means), with one for each response variable, can be computed directly as a linear transformation of Y . Provide simple expressions for the transformation. Also provide an explicit form for the joint distribution of the set of p averages.

16.4.3 For any full-rank model, $\hat{B} \sim \mathcal{N}_{q,p}[\mathbf{B}, (\mathbf{X}'\mathbf{X})^{-1}, \Sigma]$ (here $q = G$).

(a) Provide a simple expression for $(\mathbf{X}'\mathbf{X})^{-1}$.

(b) With $C = \mathbf{1}'_G/G$ and $U = \mathbf{I}_p$, provide an explicit form for the distribution of $\hat{\Theta}$ (mean of the group means, one per variable).

(c) Testing the hypothesis defined by $C = \mathbf{1}'_G/G$ and $U = \mathbf{I}_p$ involves $S_h = \hat{\Theta}' M^{-1} \hat{\Theta}$. Provide an explicit form for M^{-1} . What relevance, if any, does the term “harmonic mean” have?

(d) Assuming the alternative hypothesis holds, name the distribution and provide explicit forms for the parameters of the distribution of S_h .

16.5 Assume $y_{ig} \sim \mathcal{N}_1(\mu_g, \sigma_g^2)$ for constant $0 < \sigma_g^2 < \infty$, constant finite μ_g , $i \in \{1, \dots, n\}$, and $g \in \{1, 2\}$, with y_{ig} independent of $y_{i'g'}$ unless $i = i'$ and $g = g'$. The exercise concerns the model

$$\mathbf{y}_s = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}_s$$

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_n \end{bmatrix} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \end{bmatrix},$$

with $\mathbf{y}'_g = [y_{1g} \cdots y_{ng}]$ and testing hypothesis 1, $H_0 : \mu_1 = \mu_2$ versus $H_A : \mu_1 \neq \mu_2$.

16.5.1 Fully specify the distribution of \mathbf{e}_s .

16.5.2 Specify the likelihood function for \mathbf{y}_s .

16.5.3 Describe maximum likelihood estimators for $\{\mu_1, \mu_2, \sigma_1^2, \sigma_2^2\}$. You do not need to verify the validity of solution as a maximum (to save time in the exercises).

16.5.4 Briefly indicate why the likelihood ratio test statistic for hypothesis 1 is not a one-to-one function of a t (or F) random variable, such as occurs for the usual independent groups t .

16.5.5 Alternately, consider testing hypothesis 2, $H_0 : (\mu_1 = \mu_2) \cap (\sigma_1^2 = \sigma_2^2)$, versus $H_A : (\mu_1 \neq \mu_2) \cup (\sigma_1^2 \neq \sigma_2^2)$. Give an explicit form for the maximum log likelihood under H_0 .

Hint: Take advantage of special cases of results covered in the text. Doing so will greatly reduce the work for the exercise.

16.6 (*Optional, noncredit.*) Assume $[\text{row}_i(\mathbf{Y}_g)]' = \mathbf{Y}'_{ig} \sim \mathcal{N}_p(\boldsymbol{\mu}_g, \boldsymbol{\Sigma}_g)$ for constant $0 < |\boldsymbol{\Sigma}_g| < \infty$, constant finite $\boldsymbol{\mu}_g$, $i \in \{1, \dots, n\}$, and $g \in \{1, 2\}$, with \mathbf{Y}_{ig} independent of $\mathbf{Y}_{i'g'}$ unless $i = i'$ and $g = g'$. This question concerns the model

$$\mathbf{Y}_s = \mathbf{X}\mathbf{B} + \mathbf{E}_s$$

$$\begin{bmatrix} \mathbf{Y}_1 \\ \mathbf{Y}_2 \end{bmatrix} = \begin{bmatrix} \mathbf{1}_n & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_n \end{bmatrix} \begin{bmatrix} \boldsymbol{\mu}'_1 \\ \boldsymbol{\mu}'_2 \end{bmatrix} + \begin{bmatrix} \mathbf{E}_1 \\ \mathbf{E}_2 \end{bmatrix}$$

and testing hypothesis 1, $H_0 : \boldsymbol{\mu}_1 = \boldsymbol{\mu}_2$, versus $H_A : \boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2$.

16.6.1 Fully specify the distribution of \mathbf{E}_s .

16.6.2 Specify the likelihood function for \mathbf{Y}_s .

16.6.3 Describe maximum likelihood estimators for $\{\boldsymbol{\mu}_1, \boldsymbol{\mu}_2, \boldsymbol{\Sigma}_1, \boldsymbol{\Sigma}_2\}$. You do not need to verify the validity of solution as a maximum (to save time in the exercises).

16.6.4 Briefly indicate why the likelihood ratio test statistic for hypothesis 1 is *not* a one-to-one function of a t (or F) random variable, such as occurs for the usual independent groups Hotelling T^2 .

16.6.5 Alternately, consider testing hypothesis 2, $H_0 : (\boldsymbol{\mu}_1 = \boldsymbol{\mu}_2) \cap (\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2)$, versus $H_A : (\boldsymbol{\mu}_1 \neq \boldsymbol{\mu}_2) \cup (\boldsymbol{\Sigma}_1 \neq \boldsymbol{\Sigma}_2)$. Give an explicit form for the maximum log likelihood under H_0 .

CHAPTER 17

Tests for Generalizations of Multivariate Linear Models

17.1 MOTIVATION

Generalizations of the multivariate linear model typically allow straightforward estimation. In contrast, accurate inference in small samples often proves difficult.

Example 17.1 A $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ for blood sugar measured once per day at noon on Monday–Friday has $N \times 5$ \mathbf{Y} . The scientists wish to assess how well the amount of sugars consumed at breakfast (measured once per day, Monday–Friday) predicts blood sugar levels. A model with only linear effects is

$$\mathbf{Y} = \mathbf{X}\mathbf{B}_1 + \mathbf{E}_1 \quad (17.1)$$

$$[\mathbf{y}_M \mathbf{y}_{Tu} \mathbf{y}_W \mathbf{y}_{Th} \mathbf{y}_F] = [\mathbf{1}_N \mathbf{x}_M \mathbf{x}_{Tu} \mathbf{x}_W \mathbf{x}_{Th} \mathbf{x}_F] \begin{bmatrix} \beta_{0,M} & \beta_{0,Tu} & \beta_{0,W} & \beta_{0,Th} & \beta_{0,F} \\ \beta_{1,M} & \beta_{1,Tu} & \beta_{1,W} & \beta_{1,Th} & \beta_{1,F} \\ \beta_{2,M} & \beta_{2,Tu} & \beta_{2,W} & \beta_{2,Th} & \beta_{2,F} \\ \beta_{3,M} & \beta_{3,Tu} & \beta_{3,W} & \beta_{3,Th} & \beta_{3,F} \\ \beta_{4,M} & \beta_{4,Tu} & \beta_{4,W} & \beta_{4,Th} & \beta_{4,F} \\ \beta_{5,M} & \beta_{5,Tu} & \beta_{5,W} & \beta_{4,Th} & \beta_{5,F} \end{bmatrix} + \mathbf{E}_1.$$

Unfortunately, the model uses sugar consumed on Friday to help predict blood sugar on Monday. A more reasonable model constrains all illogical coefficients to be zero, with nonzero coefficients only for predictors on or before the day:

$$\mathbf{Y} = \mathbf{X}\mathbf{B}_2 + \mathbf{E}_2 \quad (17.2)$$

$$[\mathbf{y}_M \mathbf{y}_{Tu} \mathbf{y}_W \mathbf{y}_{Th} \mathbf{y}_F] = [\mathbf{1}_N \mathbf{x}_M \mathbf{x}_{Tu} \mathbf{x}_W \mathbf{x}_{Th} \mathbf{x}_F] \begin{bmatrix} \beta_{0,M} & \beta_{0,Tu} & \beta_{0,W} & \beta_{0,Th} & \beta_{0,F} \\ \beta_{1,M} & \beta_{1,Tu} & \beta_{1,W} & \beta_{1,Th} & \beta_{1,F} \\ 0 & \beta_{2,Tu} & \beta_{2,W} & \beta_{2,Th} & \beta_{2,F} \\ 0 & 0 & \beta_{3,W} & \beta_{3,Th} & \beta_{3,F} \\ 0 & 0 & 0 & \beta_{4,Th} & \beta_{4,F} \\ 0 & 0 & 0 & 0 & \beta_{5,F} \end{bmatrix} + \mathbf{E}_2.$$

A simpler model assumes that only the covariate from the same day is a predictor:

$$Y = XB_3 + E_3 \quad (17.3)$$

$$\begin{bmatrix} \mathbf{y}_M & \mathbf{y}_{Tu} & \mathbf{y}_W & \mathbf{y}_{Th} & \mathbf{y}_F \end{bmatrix} = \begin{bmatrix} \mathbf{1}_N & \mathbf{x}_M & \mathbf{x}_{Tu} & \mathbf{x}_W & \mathbf{x}_{Th} & \mathbf{x}_F \end{bmatrix} \begin{bmatrix} \beta_{0,M} & \beta_{0,Tu} & \beta_{0,W} & \beta_{0,Th} & \beta_{0,F} \\ \beta_{1,M} & 0 & 0 & 0 & 0 \\ 0 & \beta_{2,Tu} & 0 & 0 & 0 \\ 0 & 0 & \beta_{3,W} & 0 & 0 \\ 0 & 0 & 0 & \beta_{4,Th} & 0 \\ 0 & 0 & 0 & 0 & \beta_{5,F} \end{bmatrix} + \mathbf{E}_2.$$

All three models have a simple Gaussian likelihood. Derivatives give equations that can be solved iteratively for estimators of \mathbf{B}_j and Σ , and the likelihood ratio test is simple to compute. However, any pair of $\{\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3\}$ differ due to *nonlinear* constraints, which disallows applying standard linear hypothesis theory.

17.2 DOUBLY MULTIVARIATE MODELS

As discussed in Chapter 4, some doubly multivariate settings may be analyzed with combinations of multivariate models. The approach has the advantage of exact properties (often a Bonferroni correction must control overall test size).

Timm (Section 6.7, 2002) reviewed doubly multivariate models based on a direct-product covariance assumption. The work of Boik (1988, 1991) deserves special attention because his methods appear to control test size in small samples.

17.3 MISSING RESPONSES IN MULTIVARIATE LINEAR MODELS

Section 3.12 includes a brief discussion of missing data in the multivariate linear model. For data missing at random, the methods of Cattellier and Muller (2000) very nearly control test size, even in small samples ($N = 12$). In contrast, current mixed model tests may dramatically inflate test size. Section A.2 in the Appendix contains a description of free software (which may be downloaded from the Web) to implement the methods.

17.4 EXACT AND APPROXIMATE WEIGHTED LEAST SQUARES

Knowing exact weights allows transforming the linear model to a new model which exactly follows the usual assumptions. Some commercial software incorporates options to simplify the process.

Approximate and iterated approximate weighted least squares typically leads to very optimistic tests and inference, at least in small samples. Some form of corrected or test should be used, if available. Estimation is typically well behaved.

17.5 SEEMINGLY UNRELATED REGRESSIONS

Seemingly unrelated regression models (Srivastava and Giles, 1987) also allow accurate estimation. However, the specific nature of the model has allowed statisticians to develop approximate tests which work reasonably well in small samples (Rocke, 1989). Also, Timm (2002, Section 5.14) gave a useful review of testing, including some newer work.

17.6 GROWTH CURVE MODELS (GMANOVA)

From the perspective of today, early formulations of growth curve models may be transformed exactly into special cases of the general linear multivariate model. Often a GCM expresses repeated measures as a polynomial function of time. A quadratic model uses only the average, linear, and quadratic trends and hence $p_T = 3$ coefficients. The approach corresponds to restricting attention to a reduced model defined by a transformation. For a $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with $p = 4$ equally spaced repeated measures,

$$U_T = \begin{bmatrix} 1 & -3 & 1 \\ 1 & -1 & -1 \\ 1 & 1 & -1 \\ 1 & 3 & 1 \end{bmatrix} \begin{bmatrix} 4 & 0 & 0 \\ 0 & 20 & 0 \\ 0 & 0 & 4 \end{bmatrix}^{-1/2} \tag{17.4}$$

defines an $N \times 3$ matrix of lower order trend scores,

$$Y_T = Y U_T. \tag{17.5}$$

Transforming the original model gives

$$\begin{aligned} Y U_T &= X B U_T + E U_T \\ Y_T &= X B_T + E_T. \end{aligned} \tag{17.6}$$

The theory of the multivariate linear model applies exactly to the reduced model.

The $N \times 1$ set of cubic trend scores has been ignored in the reduced model. Later formulations introduced the use of the ignored high-order trends as covariates to increase precision of the estimators. Unfortunately, the process makes hypothesis tests and confidence intervals optimistic (Berger, 1986).

17.7 TESTING HYPOTHESES IN THE GCM

The GCM is a multivariate GLM with Gaussian errors and $p - q$ linear restrictions on B . The unrestricted model is $GLM_{N,p,q}(Y_i; X_i B_*, \Sigma)$, while the restrictions give $GLM_{N,p,m}(Y_i; X_i B_* | B_* = B_1 T_1, \Sigma)$ with B_* ($q \times p$). Equivalently

$$E(Y) = X B_* \tag{17.7}$$

$$\mathcal{V}[\text{vec}(Y')] = I \otimes \Sigma. \tag{17.8}$$

The general GCM is obtained from the GLM by choosing $B_* = B_1 T_1$ ($q \times m \times p$) or, equivalently,

$$B_* = [B_1 \ B_2] \begin{bmatrix} T_1 \\ T_2 \end{bmatrix}, \tag{17.9}$$

with $B_2 = \mathbf{0}$. Here we assume T ($p \times p$) has full rank and $T^{-1} = G = [G_1 \ G_2]$, with G_1 ($p \times m$) and G_2 [$p \times (p - m)$]. The relationship between B_* and B can be written as $B_* [G_1 \ G_2] = [B_1 \ B_2]$, with B_1 ($q \times m$) and B_2 [$q \times (p - m)$]. Since $B_2 = \mathbf{0}$, the equation makes it clear an equivalent notation is $GLM_{N,p,q}(Y_i; X_i B_* | B_* G_2 = \mathbf{0}, \Sigma)$ with G_2 of dimension $p \times (p - m)$. The matrix T has specified structure, such as a polynomial structure, and the restriction imposes the structure on B_* ($q \times p$) in the form of $p - m$ linearly independent constraints. With T based on polynomials, the restriction indicates the elements of B_* are functionally related to one another as specified by the polynomial.

Example 17.2 If the children's weights in Example 13.1 had been recorded monthly from the ages of 5 to 7 years, then B_* would be 1×25 . If growth is truly linear during the time of observation, then only two parameters (intercept and slope) are needed to characterize the growth curve. The 25 elements of B_* are the mean weights over 25 months. Since the elements are functionally related (they form a straight line), the 25 elements (means) can be condensed into 2 columns:

$$B_* = [\beta_0 \ \beta_1 \ 0 \ \cdots \ 0] T. \tag{17.10}$$

Therefore it is appropriate to impose the restriction $B_* G_2 = \mathbf{0}$, with G_2 (25×23).

Theorem 17.1 The $GCM_{N,p,q,m}(Y_i; X_i B_1 T_1, \Sigma)$ with B and T defined in terms of an orthogonal or orthonormal polynomial allows testing $H_0 : C B_1 U = \mathbf{0}$ ($a \times b$) against the general alternative. Under $GLM_{N,p,q}(Y_i; X_i B_* | B_* G_2 = \mathbf{0}, \Sigma)$ with G_2 of dimension $p \times (p - m)$, an equivalent hypothesis is $H_0 : C B_* G_1 U = \mathbf{0}$.

Proof. With $T^{-1} = G = [G_1 \ G_2]$ we have $B_* G_1 = B_1$. □

A more general version of the linear restrictions on B_* involves an arbitrary constant matrix D , with $B_* = B T + D$. No new problems arise. The form leads to an ordinary GCM having expectation of the form $E(Y - X D) = X B T$.

17.8 CONFIDENCE BANDS FOR GROWTH CURVES

For the polynomial $GCM_{N,p,q,m}(Y_i; X_i B T, \Sigma)$ with $m = 2$, Stewart (1987) developed a finite-interval confidence band in terms of a type III bivariate t distribution, in a manner analogous to Bowden's (1970) univariate results. The procedure generalizes to $m > 2$ in terms of type III multivariate t distributions.

Tests for Linear Mixed Models

18.1 OVERVIEW

We introduced the notation $\text{LMM}_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\boldsymbol{\beta}, \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}_i' + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)]$ in Chapter 5 to indicate a general linear mixed model, with model equation $\mathbf{y}_i = \mathbf{X}_i\boldsymbol{\beta} + \mathbf{Z}_i\mathbf{d}_i + \mathbf{e}_i$ for one independent sampling unit (ISU). In the present chapter, we always add the Gaussian assumption, which has three parts: $\mathbf{e}_i \sim \mathcal{N}_p[\mathbf{0}, \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)]$, $\mathbf{d}_i \sim \mathcal{N}_q[\mathbf{0}, \boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)]$, and $\mathbf{e}_i \perp\!\!\!\perp \mathbf{d}_i$. The model implies $E(\mathbf{y}_i) = \mathbf{X}_i\boldsymbol{\beta}$ and $\mathcal{V}(\mathbf{y}_i) = \mathbf{Z}_i\boldsymbol{\Sigma}_{d_i}(\boldsymbol{\tau}_d)\mathbf{Z}_i' + \boldsymbol{\Sigma}_{e_i}(\boldsymbol{\tau}_e)$. Hence $\boldsymbol{\beta}$ determines the expected values, the first moments of the observations. In turn, $\boldsymbol{\tau}_d$ determines the covariance matrix, the second moments, of the unobserved (latent) error components varying within the ISU due to observed characteristics in \mathbf{Z}_i . Similarly, $\boldsymbol{\tau}_e$ determines the covariance matrix of the unobserved (latent) error components varying within the ISU due to observation. The Gaussian structure ensures parameters $\boldsymbol{\beta}$ and $\boldsymbol{\tau}' = [\boldsymbol{\tau}_d' \boldsymbol{\tau}_e']$ fully determine the distributions of all model statistics.

The presence of just mean and covariance parameters leads to three possible types of inferences. Most often, data analysts seek to draw inferences about functions of $\boldsymbol{\beta}$, the parameters determining the means (the fixed effects). Less often, data analysts seek to draw inferences about functions of $\boldsymbol{\tau}$, the covariance parameters, especially $\boldsymbol{\tau}_d$ (the random effects parameters). Even less often, data analysts seek to draw inferences about functions directly involving both $\boldsymbol{\beta}$ and $\boldsymbol{\tau}$.

The univariate and multivariate models allow computing exact or nearly exact likelihood-based inferences as an adjunct to the noniterative calculation of estimates of a single model. In contrast, the mixed model requires iterative calculation of estimates for two or more models. The desire for speed and convenience has led to a heavy emphasis on approximate tests based on single-model fits. The inexorable increase in computer speed has made the computer time needed to fit an additional model much less consequential. However, the control language for contemporary software makes comparing two model fits awkward and hence creates a barrier to conducting tests based on two model fits.

At the present time, confidence intervals and hypothesis tests based on a single model fit dominate practice. Two related reasons reinforce the habit. First, popular contemporary commercial software does not directly implement any

inference based on fitting two or more models. Second, invoking software twice to fit two distinct models takes extra time and care for the data analyst to implement correctly.

Vonesh and Chinchilli (1997), Verbeke and Molenberghs (2000), and Demidenko (2004) provided excellent book-length treatments of mixed models. Khuri, Mathew, and Sinha (1998) focused specifically on tests in mixed models.

18.2 ESTIMABILITY OF $\theta = C\beta$

For an estimable secondary parameter $\theta = C\beta$, we consider testing

$$H_0 : \theta = \theta_0. \quad (18.1)$$

In parallel to extending the concept of estimable θ to mixed models, the concept of a testable hypothesis extends to the mixed model. In particular, we require estimable θ and full (row) rank of C . Results about testability center on guaranteeing existence of a valid test. Such statements involve population parameter and design matrix properties and hence apply to mixed models.

18.3 LIKELIHOOD RATIO TESTS OF $C\beta$

Hypotheses concerning $\theta = C\beta$ can be tested with a likelihood ratio test. However the distribution theory is not well developed. The process involves finding estimates for two models, with one of them false. The iterative calculations for the false model may be unstable. Although theoretically appealing, the method is less used than methods based on fitting a single model.

Example 18.1 The likelihood ratio approach allows testing $H_0 : C\beta = \theta_0$ in the $GLM_{N,q}(y_i; \mathbf{X}_i\beta, \sigma^2)$ with (i.i.d.) Gaussian errors, as a special case of the mixed model. With $r = \text{rank}(\mathbf{X})$, the exact distribution of the likelihood ratio test statistic can be expressed in terms of

$$F = \frac{(\hat{\theta} - \theta_0)' \mathbf{M}^{-1} (\hat{\theta} - \theta_0) / a}{\hat{\sigma}^2} = \frac{SSH/a}{SSE/(N-r)} \sim F(a, N-r). \quad (18.2)$$

The mixed model likelihood ratio test approximation corresponds to saying $(Fa) \sim \chi^2(a)$, which is correct only asymptotically. With $SSE/\sigma^2 = \hat{\sigma}^2/\sigma^2 \sim \chi^2(N-r)$, $\hat{\sigma}^2 \rightarrow \sigma^2$ as $N \rightarrow \infty$. For any N , under the null $SSH/\sigma^2 \sim \chi^2(a)$. The denominator degrees-of-freedom parameter of the F distribution accounts for having estimated σ^2 in a finite sample.

It is typically straightforward, and always possible, to code a linear model design matrix so that a hypothesis test corresponds to deleting one or more variables from the model. Doing so allows computing the likelihood ratio test

statistic as the difference between the log likelihood values for the full and reduced model. The asymptotic approximation is $-2\ln[L(\hat{\theta}|C\beta = \theta_0)/L(\hat{\theta})] \sim \chi^2(a)$. Unfortunately, more accurate approximations are available only for some special cases. Zucker, Lieberman, and Manor (2000) described a second-order approximation for a particular class of models and also demonstrated the substantial improvement in accuracy that the approach provides.

18.4 LIKELIHOOD RATIO TESTS INVOLVING τ

Many linear and some nonlinear hypotheses concerning τ alone or $\gamma' = [\beta' \tau']$ can be tested with the likelihood ratio approach. As always, the hypothesis of interest must correspond to comparing a full and constrained model with the constrained model nested within the full. Maximum likelihood estimates are used. Some special case tests for some models have known distributions (Demidenko, 2004). Otherwise only large-sample results are available.

18.5 TEST SIZE OF WALD-TYPE TESTS OF β USING REML

Among the widely available approximate F methods based on fitting a single model, the Kenward and Roger (1997) technique appears to provide the most accurate test size. However, substantial room for improvement in small-sample performance remains. Simulations of Park, Park, and Davis (2001) as well as simulations of Schaalje, McBride, and Fellingham (2003) support the conclusion.

Table 18.1 Estimated Test Size, *Target* = 0.04, for UNIREP (Std. Err. < 0.0004) and Mixed (Std. Err. < 0.003)

<i>N</i>	ϵ	UNIREP		Mixed		
		GG ¹	HF ²	Resid ³	Satter ⁴	K/R ⁵
10	0.28	0.042	0.045	0.254	0.138	0.114
	0.51	0.039	0.052	0.263	0.137	0.081
	1.00	0.021	0.052	0.263	0.144	0.043
20	0.28	0.041	0.042	0.116	0.077	0.040
	0.51	0.040	0.046	0.115	0.072	0.036
	1.00	0.029	0.038	0.116	0.075	0.038
40	0.28	0.040	0.041	0.075	0.054	0.040
	0.51	0.041	0.043	0.076	0.061	0.045
	1.00	0.034	0.039	0.074	0.056	0.040

¹Geisser Greenhouse, ²Huynh-Feldt

³Residual, ⁴Satterthwaite, ⁵Kenward Roger

Muller, Edwards, Simpson and Taylor (2006) illustrated some limitations of the approach with a simulation summarized in Table 18.1. Simulated Gaussian data

were generated which met all assumptions of a $GLM_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \Sigma)$ with $p = 9$ and $q = 1$, with $\mathbf{X} = \mathbf{1}_N$, a grand mean model. No missing or mistimed data were allowed. The 9 responses reflected a 3×3 within-subject factorial design, with factor names Clip and Region. The study design used $\alpha = 0.04$, not 0.05. For the UNIREP test of Clip \times Region interaction, $H_0 : \mathbf{C} \mathbf{B} \mathbf{U} = \Theta_0$ used $\mathbf{C} = \mathbf{1}$, $\Theta_0 = \mathbf{0}$, and

$$\begin{aligned}
 U_{cr} &= \mathbf{T}_c \otimes \mathbf{T}_r \\
 &= \begin{bmatrix} -4/\sqrt{42} & 2/\sqrt{14} \\ -1/\sqrt{42} & -3/\sqrt{14} \\ 5/\sqrt{42} & 1/\sqrt{14} \end{bmatrix} \otimes \begin{bmatrix} -1/\sqrt{2} & 1/\sqrt{6} \\ 0 & -2/\sqrt{6} \\ 1/\sqrt{2} & 1/\sqrt{6} \end{bmatrix}. \quad (18.3)
 \end{aligned}$$

Columns of \mathbf{T}_c contain the linear and quadratic orthonormal trends for $\log_2(\text{Clip}) \in \{1, 2, 4\}$, and columns of \mathbf{T}_r contain the orthonormal contrasts of linear and quadratic trends for $\log_2(\text{Region}) \in \{1, 3, 5\}$. Four eigenvalues sets, $\lambda'_1 \approx [0.4796 \ 0.0100 \ 0.0100 \ 0.0100]$, $\lambda'_2 \approx [0.3455 \ 0.0612 \ 0.0556 \ 0.0472]$, $\lambda'_3 \approx [0.2355 \ 0.1712 \ 0.0556 \ 0.0472]$, and $\lambda'_4 \approx [0.1274 \ 0.1274 \ 0.1274 \ 0.1274]$, were used. Corresponding values of ϵ are approximately 0.28, 0.51, 0.72, and 1. Having $\epsilon = 1$ corresponds to having the underlying repeated measures being compound symmetric. The conditions are the same as conditions 5–8 in Table III in Coffey and Muller (2003). Given $\Sigma_* = \text{Dg}(\lambda_j)$ for $j \in \{1, 2, 3, 4\}$, $\Sigma = U_{cr} \Sigma_* U'_{cr}$. A total of 500,000 replications were tabulated for each condition for the UNIREP tests (Geisser-Greenhouse and Huynh-Feldt). A total of 5000 replications were tabulated for the mixed tests based on fitting a single model with SAS PROC MIXED[®] and using an F approximation (residual sum of squares, Satterthwaite, and Kenward-Roger approximations for denominator degrees of freedom). The “unstructured” covariance option on the repeated statement was always used. In addition to the type I error rate inflation, the mixed model approach also failed to converge in a fraction of the cases. A number of adjustments to the inputs controlling the estimation algorithms greatly reduced but did not completely eliminate the problem. Each set of data which led to convergence failure was subsequently analyzed with the UNIREP approach, which always gave well-behaved and reasonable estimates for the data observed. As noted earlier, no missing or mistimed data were present.

In interpreting the results, it helps to keep in mind the distinction between number of independent sampling units (ISUs), with $N \in \{10, 20, 40\}$ in the simulations, and number of observations, with $n = 9N \in \{90, 180, 360\}$ in the simulations. The residual sums-of-squares approximation inflates test size the most. It implicitly uses weighted least squares estimators and makes no adjustment for having estimated the covariance parameters. Consequently, denominator degrees of freedom are based on $n - \text{rank}(\mathbf{X}_s)$, which coincides with the degrees of freedom for the uncorrected UNIREP test. The corrected tests use reduced numerator and denominator degrees of freedom.

18.6 USING WALD-TYPE TESTS OF β WITH REML

We next describe the steps needed to compute the Kenward and Roger (1997) approximation for inference about elements of β . The approach assumes $\Sigma_i(\tau)$ has been correctly specified and does not change. Except for notation differences, our presentation closely follows the original.

Using notation introduced in Chapter 5, we consider $LMM_{N,p,q,m}[\mathbf{y}_i; \mathbf{X}_i\beta, \mathbf{Z}_i\Sigma_{di}(\tau_d)\mathbf{Z}'_i + \Sigma_{ei}(\tau_e)]$ with Gaussian errors. The model for all of the data stacked by participant (ISU) may be written

$$\begin{aligned} \mathbf{y}_s &= \mathbf{X}_s\beta + \mathbf{Z}_s\mathbf{d}_s + \mathbf{e}_s \\ \begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \\ \vdots \\ \mathbf{y}_N \end{bmatrix} &= \begin{bmatrix} \mathbf{X}_1 \\ \mathbf{X}_2 \\ \vdots \\ \mathbf{X}_N \end{bmatrix} \beta + \begin{bmatrix} \mathbf{Z}_1 & 0 & \cdots & 0 \\ 0 & \mathbf{Z}_2 & & \vdots \\ \vdots & & \ddots & 0 \\ 0 & \cdots & 0 & \mathbf{Z}_N \end{bmatrix} \begin{bmatrix} \mathbf{d}_1 \\ \mathbf{d}_2 \\ \vdots \\ \mathbf{d}_N \end{bmatrix} + \begin{bmatrix} \mathbf{e}_1 \\ \mathbf{e}_2 \\ \vdots \\ \mathbf{e}_N \end{bmatrix}. \end{aligned} \tag{18.4}$$

For an $a \times 1$ testable secondary parameter $\theta = \mathbf{C}\beta$, necessarily with $\text{rank}(\mathbf{C}) = a \leq q$, we consider testing the hypothesis

$$H_0 : \theta = \theta_0 \tag{18.5}$$

against the general alternative. Here $E(\mathbf{y}_i) = \mathbf{X}_i\beta$ implies $E(\mathbf{y}_s) = \mathbf{X}_s\beta$, while $\Sigma_i(\tau) = \mathbf{Z}_i\Sigma_{di}(\tau_d)\mathbf{Z}'_i + \Sigma_{ei}(\tau_e)$ implies

$$\mathcal{V}(\mathbf{y}_s) = \bigoplus_{i=1}^N \Sigma_i(\tau) = \Sigma_s(\tau). \tag{18.6}$$

The unbiased REML estimator of β ($q \times 1$) is

$$\hat{\beta} = [\mathbf{X}'_s \Sigma_s^{-1}(\hat{\tau}) \mathbf{X}_s]^{-1} \mathbf{X}'_s \Sigma_s^{-1}(\hat{\tau}) \mathbf{y}_s. \tag{18.7}$$

The asymptotic covariance matrix of $\hat{\beta}$ is

$$\mathcal{V}_a(\hat{\beta}) = \Phi(\tau) = [\mathbf{X}'_s \Sigma_s^{-1}(\tau) \mathbf{X}_s]^{-1}. \tag{18.8}$$

More precisely, with s elements in τ ,

$$[\Phi(\tau)]^{-1/2} (\hat{\beta} - \beta) \overset{\text{asy}}{\approx} \mathcal{N}_s(\mathbf{0}, \mathbf{I}). \tag{18.9}$$

As always, the distinction between the number of ISUs, N , and the number of observations, $n = \sum_{i=1}^N p_i$, must be treated carefully. Fixing N and having $p_i \rightarrow \infty$ implies $n \rightarrow \infty$ but does not guarantee convergence of $\mathcal{V}(\hat{\beta})$, except with side conditions to greatly simplify the covariance structure.

The following asymptotically correct estimator has often been used for evaluating the precision of $\hat{\beta}$ and for testing hypotheses:

$$\begin{aligned} \tilde{V}_1(\hat{\beta}) &= \Phi(\hat{\tau}) = [X'_s \Sigma_s^{-1}(\hat{\tau}) X_s]^{-1} \\ &= \tilde{\Phi}_1. \end{aligned} \tag{18.10}$$

The estimator treats $\hat{\tau}$ in equation (18.7) as if it were a fixed constant. Doing so corresponds to assuming $\mathcal{V}(\hat{\tau}) \approx \mathbf{0}$, which is not true for finite sample sizes. In turn, $\Phi(\hat{\tau}) = \tilde{\Phi}_1$ is a biased estimator of $\Phi(\tau)$, which is an approximation of $\mathcal{V}(\hat{\beta})$. Consequently, when N is small, $\tilde{V}_1(\hat{\beta})$ underestimates $\mathcal{V}(\hat{\beta})$ due to a combination of bias and approximation error.

Kackar and Harville (1984) used a first-order Taylor series expansion to better approximate $\mathcal{V}(\hat{\beta})$ in small sample sizes. The expansion depends on

$$P_j = X'_s \left\{ \frac{\partial \Sigma_s^{-1}(\tau)}{\partial \tau_j} \right\} X_s \tag{18.11}$$

and

$$Q_{jk} = X'_s \left\{ \frac{\partial \Sigma_s^{-1}(\tau)}{\partial \tau_j} \right\} \Sigma_s(\tau) \left\{ \frac{\partial \Sigma_s^{-1}(\tau)}{\partial \tau_k} \right\} X_s. \tag{18.12}$$

The expansion leads to the approximation

$$\mathcal{V}(\hat{\beta}) \approx \Phi(\tau) + \Phi(\tau) \left\{ \sum_{j=1}^t \sum_{k=1}^t \langle \mathcal{V}(\hat{\tau}) \rangle_{jk} [Q_{jk} - P_j \Phi(\tau) P_k] \right\} \Phi(\tau). \tag{18.13}$$

An estimator of $\mathcal{V}(\hat{\beta})$ may be computed in three steps: (1) evaluate P_j and Q_{jk} at $\hat{\tau}$ to give \hat{P}_j and \hat{Q}_{jk} ; (2) replace $\Phi(\tau)$ by $\tilde{\Phi}_1 = \Phi(\hat{\tau})$ from equation 18.10; and (3) replace $\mathcal{V}(\hat{\tau})$ with an estimator, $\tilde{\mathcal{V}}(\hat{\tau})$. The value of $\tilde{\mathcal{V}}(\hat{\tau})$ may be computed either as the inverse of the τ submatrix of the expected information matrix or as the τ submatrix of the inverse of the observed information matrix. In turn,

$$\tilde{V}_2(\hat{\beta}) = \tilde{\Phi}_1 + \tilde{\Phi}_1 \left[\sum_{j=1}^t \sum_{k=1}^t \langle \tilde{\mathcal{V}}(\hat{\tau}) \rangle_{jk} (\hat{Q}_{jk} - \hat{P}_j \tilde{\Phi}_1 \hat{P}_k) \right] \tilde{\Phi}_1. \tag{18.14}$$

Kenward and Roger (1997) used the Kackar and Harville result to devise a further improved estimator of $\mathcal{V}(\hat{\beta})$. The estimator adds an adjustment in estimating $\Phi(\hat{\tau})$ to compensate for finite N [and $\mathcal{V}(\hat{\tau}) \neq \mathbf{0}$]. With

$$R_{jk} = X'_s \Sigma_s^{-1}(\tau) \left\{ \frac{\partial^{(2)} \Sigma_s^{-1}(\tau)}{\partial \tau_j \partial \tau_k} \right\} [\Sigma_s^{-1}(\tau)] X_s / 4 \tag{18.15}$$

and $\widehat{\mathbf{R}}_{jk}$ the value of \mathbf{R}_{jk} evaluated at $\widehat{\boldsymbol{\tau}}$, the adjusted estimator is

$$\begin{aligned} \widetilde{\mathcal{V}}_3(\widehat{\boldsymbol{\beta}}) &= \widetilde{\boldsymbol{\Phi}}_1 + 2\widetilde{\boldsymbol{\Phi}}_1 \left[\sum_{j=1}^t \sum_{k=1}^t \langle \widetilde{\mathcal{V}}(\widehat{\boldsymbol{\tau}}) \rangle_{jk} \left(\widehat{\mathbf{Q}}_{jk} - \widehat{\mathbf{P}}_j \widetilde{\boldsymbol{\Phi}}_1 \widehat{\mathbf{P}}_k - \widehat{\mathbf{R}}_{jk} \right) \right] \widetilde{\boldsymbol{\Phi}}_1. \\ &= \widetilde{\boldsymbol{\Phi}}_2. \end{aligned} \tag{18.16}$$

Kenward and Roger found, in limited simulations, that replacing the expected information matrix by the observed information matrices had no discernible effect on accuracy of tests and standard errors. However, the missing data mechanisms in their examples and simulations were arguably MCAR. More generally, many kinds of MAR mechanisms can require using the observed information matrix to avoid bias (Verbeke and Molenberghs, 2000, Chapter 21).

Kenward and Roger (1997) provided an F approximation for hypothesis tests. Computing the test statistic, F_{KR} , requires a denominator degrees of freedom and a scaling constant λ . Both depend on

$$\mathbf{T} = \mathbf{C}'(\mathbf{C}\boldsymbol{\Phi}^{-1}\mathbf{C}')^{-1}\mathbf{C}, \tag{18.17}$$

$$d_1 = \sum_{j=1}^t \sum_{k=1}^t \langle \mathcal{V}(\widehat{\boldsymbol{\tau}}) \rangle_{jk} \text{tr}(\mathbf{T}\boldsymbol{\Phi}\mathbf{P}_j\boldsymbol{\Phi}) \text{tr}(\mathbf{T}\boldsymbol{\Phi}\mathbf{P}_k\boldsymbol{\Phi}), \tag{18.18}$$

$$d_2 = \sum_{j=1}^t \sum_{k=1}^t \langle \mathcal{V}(\widehat{\boldsymbol{\tau}}) \rangle_{jk} \text{tr}(\mathbf{T}\boldsymbol{\Phi}\mathbf{P}_j\boldsymbol{\Phi}\mathbf{T}\boldsymbol{\Phi}\mathbf{P}_k\boldsymbol{\Phi}), \tag{18.19}$$

$g = [(a+1)d_1 - (a+4)d_2] / [(a+2)d_2]$, $c_1 = g / [3a+2(1-g)]$, $c_2 = (a-g) / [3a+2(1-g)]$, $c_3 = (a-g+2) / [3a+2(1-g)]$, $b = (d_1+6d_2) / (2a)$, $c_4 = (1-d_2/a)^{-1}$, $c_5 = 2(1+c_1b) / [a(1-c_2b)^2(1-c_3b)]$, and $\rho = c_5 / (2c_4)^2$. In turn,

$$\nu_* = 4 + (a+2) / (a\rho - 1) \tag{18.20}$$

and

$$\lambda = \nu_* / [(\nu_* - 2)c_4]. \tag{18.21}$$

The simplest approximation uses $F_m = (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)'(\mathbf{C}\widetilde{\boldsymbol{\Phi}}_1^{-1}\mathbf{C}')^{-1}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) / a$. Replacing $\widetilde{\boldsymbol{\Phi}}_1$ by $\widetilde{\boldsymbol{\Phi}}_2$ in the expression for F_m and using $\widetilde{\boldsymbol{\Phi}}_2$ to compute $\widehat{\nu}_*$ and $\widehat{\lambda}$ gives

$$\begin{aligned} F_{\text{KR}} &= \widehat{\lambda} \cdot (\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)'(\mathbf{C}\widehat{\boldsymbol{\Phi}}_2^{-1}\mathbf{C}')^{-1}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0) / a \\ &\sim F(a, \widehat{\nu}_*). \end{aligned} \tag{18.22}$$

18.7 USING WALD-TYPE TESTS OF $\{\beta, \tau\}$ WITH REML

A simpler approximation for tests of $H_0 : \theta = \theta_0$, with $\theta = C\beta$, involves using

$$F_m = (\hat{\theta} - \theta_0)'(C\tilde{\Phi}_1^{-1}C')^{-1}(\hat{\theta} - \theta_0)/a. \tag{18.23}$$

A variation on the test statistic was discussed in Helms (1992). The distribution of F_m is approximated by the $F(a, \nu_2, \omega)$ distribution in which $a = \text{rank}(C)$, $\nu_2 = N - \text{rank}([\mathbf{X}_s \ \mathbf{Z}_s])$, and $\omega = (\theta - \theta_0)'[C(\mathbf{X}'_s \Sigma_s^{-1} \mathbf{X}_s)^{-1}C']^{-1}(\theta - \theta_0)$. The appropriate critical value for testing H_0 versus the general alternative is $f_C = F_F^{-1}(1 - \alpha; a, \nu_2)$. The power of the approximately size- α test is approximately $\text{Power} = 1 - F(f_C; a, \nu_2, \omega)$. Only limited simulation results are available. As noted earlier, the Kenward-Roger approximation provides greater accuracy for inference about $\theta = C\beta$.

The same form of approximate test can be applied to tests involving both $\{\beta, \tau\}$. It seems likely that the roughness of approximation requires a substantial number of independent sampling units to provide the desired distribution. If, as defined in equation 5.9, $\Sigma_{ds} = \bigoplus_{i=1}^N \Sigma_{di}$ and $\Sigma_{es} = \bigoplus_{i=1}^N \Sigma_{ei}$, then

$$\nu \left(\begin{bmatrix} \hat{\beta} \\ \hat{d}_s - d_s \end{bmatrix} \right) = \begin{bmatrix} \mathbf{X}'_s \Sigma_{es}^{-1} \mathbf{X}_s & \mathbf{X}'_s \Sigma_{es}^{-1} \mathbf{Z}_s \\ \mathbf{Z}'_s \Sigma_{es}^{-1} \mathbf{X}_s & \Sigma_{ds}^{-1} + \mathbf{Z}'_s \Sigma_{es}^{-1} \mathbf{Z}_s \end{bmatrix}^{-1} = \mathbf{M}^{-1}. \tag{18.24}$$

For \mathbf{M} positive definite,

$$\begin{aligned} \mathbf{M}^{-1} &= \begin{bmatrix} \mathbf{M}^{(11)} & \mathbf{M}^{(12)} \\ \mathbf{M}^{(21)} & \mathbf{M}^{(22)} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{M}_{11} & \mathbf{M}_{12} \\ \mathbf{M}_{21} & \mathbf{M}_{22} \end{bmatrix}^{-1}. \end{aligned} \tag{18.25}$$

Tests of hypotheses involving d_s , such as

$$H_0 : L' \begin{bmatrix} \beta \\ d_s \end{bmatrix} = \theta_0, \tag{18.26}$$

are testable against the general alternative using test statistic

$$F_L = L' \left(\begin{bmatrix} \hat{\beta} \\ \hat{d}_s \end{bmatrix} - \theta_0 \right)' (L\hat{\mathbf{M}}^{-1}L')^{-1} \left(\begin{bmatrix} \hat{\beta} \\ \hat{d}_s \end{bmatrix} - \theta_0 \right) / \nu_1. \tag{18.27}$$

Here $\nu_1 = \text{rank}(L)$, $\nu_2 = n - \text{rank}([\mathbf{X}_s \ \mathbf{Z}_s])$, and

$$\omega = L' \left(\begin{bmatrix} \beta \\ d_s \end{bmatrix} - \theta_0 \right)' (LM^{-1}L')^{-1} \left(\begin{bmatrix} \beta \\ d_s \end{bmatrix} - \theta_0 \right). \tag{18.28}$$

The appropriate critical value for testing H_0 versus the general alternative is $f_C = F_F^{-1}(1 - \alpha, \nu_1, \nu_2)$. The power of the approximate size- α test is approximately $\text{Power} = 1 - F(f_C; \nu_1, \nu_2, \omega)$.

A Review of Multivariate and Univariate Linear Models

19.1 MATRIX GAUSSIAN AND WISHART PROPERTIES

The results in the present section give the basic distribution theory needed for parameter estimators in the multivariate and univariate linear models with Gaussian errors. Only the test statistic distributions require separate treatment. Univariate results occur as the special case with $p = 1$.

We begin by reproducing key results about Gaussian and Wishart distributions from Chapter 8 and Chapter 10. The results provide the basis for nearly all distributional results in the remainder of the chapter.

Copy of Definition 8.4 The $n \times p$ random matrix Y follows a *direct-product matrix Gaussian* distribution, typically abbreviated *matrix Gaussian* and written $Y \sim \mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$, if and only if

- (a) $\text{vec}(Y) \sim (S)\mathcal{N}_{n,p}[\text{vec}(\mathbf{M}), \Sigma \otimes \Xi]$; if and only if
- (b) $\text{vec}(Y') \sim (S)\mathcal{N}_{n,p}[\text{vec}(\mathbf{M}'), \Xi \otimes \Sigma]$; if and only if
- (c) $Y = \Psi Z \Phi' + \mathbf{M}$ with $\text{vec}(Z) \sim \mathcal{N}_{n_1 p_1}(\mathbf{0}, \mathbf{I})$ and
 - Ψ ($n \times n_1$) of rank $n_1 \geq 1$, $\Xi = \Psi \Psi'$,
 - Φ' ($p_1 \times p$) of rank $p_1 \geq 1$, $\Sigma = \Phi \Phi'$.

Writing $Y \sim S\mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$ indicates $n_1 = \text{rank}(\Psi) = \text{rank}(\Xi) < n$, or $p_1 = \text{rank}(\Phi) = \text{rank}(\Sigma) < p$, or both.

Writing $Y \sim (S)\mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$ emphasizes allowing any combination of $n_1 \leq n$ and $p_1 \leq p$.

Copy of Theorem 8.12 If $Y \sim (S)\mathcal{N}_{n,p}(\mathbf{M}, \Xi, \Sigma)$, while $A \neq \mathbf{0}$ ($n_1 \times n$), $B \neq \mathbf{0}$ ($p \times p_1$) and C ($n_1 \times p_1$) are finite constant matrices, then

$$AYB + C \sim (S)\mathcal{N}_{n_1, p_1}(AMB + C, A\Xi A', B'\Sigma B). \quad (19.1)$$

Copy of Definition 10.2 (a) If $Y \sim \mathcal{N}_{\nu,p}(\mathbf{0}, \mathbf{I}_\nu, \Sigma)$ then $Y'Y \sim \mathcal{W}_p(\nu, \Sigma)$ indicates $Y'Y$ follows a *central (integer) Wishart* distribution with (integer) $\nu > 0$ degrees of freedom.
 (b) If $Y \sim \mathcal{N}_{\nu,p}(\mathbf{M}, \mathbf{I}_\nu, \Sigma)$, then $Y'Y \sim \mathcal{W}_p(\nu, \Sigma, \mathbf{M}'\mathbf{M})$ indicates $Y'Y$ follows a *noncentral (integer) Wishart* distribution with (integer) $\nu > 0$ degrees of freedom, *shift* $\Delta = \mathbf{M}'\mathbf{M}$, and *noncentrality* $\Omega = \mathbf{M}'\mathbf{M}\Sigma^+$.
 (c) Singular Σ may be emphasized by writing $S\mathcal{W}_p(\nu, \Sigma)$ or $S\mathcal{W}_p(\nu, \Sigma, \Delta)$.
 (d) Writing $(S)\mathcal{W}_p(\nu, \Sigma)$ or $(S)\mathcal{W}_p(\nu, \Sigma, \Delta)$ indicates possibly singular Σ .

Copy of Theorem 10.4 If $S \sim (S)\mathcal{W}_p(\nu, \Sigma, \Delta)$ and T is any $p \times p_1$ constant matrix, then

$$T'ST \sim (S)\mathcal{W}_{p_1}(\nu, T'\Sigma T, T'\Sigma T). \tag{19.2}$$

Copy of Theorem 10.8 If $Y \sim (S)\mathcal{N}_{N,p}(\mathbf{M}, \mathbf{I}_N, \Sigma)$, $\text{rank}(\Sigma) = p_1 \leq p$ with $N \geq p_1$, $N \times N$ \mathbf{A} and \mathbf{B} are constants, the following hold.

- (a) $\mathbf{A}Y$ and $\mathbf{B}Y$ are independent if and only if $\mathbf{A}\mathbf{B}' = \mathbf{0}$.
- (b) If $\mathbf{A} = \mathbf{A}'$ is positive definite or positive semidefinite and $\mathbf{A}\mathbf{B}' = \mathbf{0}$, then $\mathbf{B}Y$ and $Y'\mathbf{A}Y$ are independent.
- (c) If $\mathbf{A} = \mathbf{A}'$ and $\mathbf{B} = \mathbf{B}'$ are positive definite or positive semidefinite and $\mathbf{B}\mathbf{A} = \mathbf{0}$, then $Y'\mathbf{A}Y$ and $Y'\mathbf{B}Y$ are independent.
- (d) If $\mathbf{A} = \mathbf{A}' = \mathbf{A}^2$, then $\nu = \text{rank}(\mathbf{A}) = \text{tr}(\mathbf{A})$ and

$$S = Y'\mathbf{A}Y \sim (S)\mathcal{W}_p(\nu, \Sigma, \mathbf{M}'\mathbf{A}\mathbf{M}). \tag{19.3}$$

Copy of Corollary 10.8.1 If $p = 1$ and $\sigma^2 > 0$, then $\mathbf{y} \sim \mathcal{N}_{n,1}(\boldsymbol{\mu}, \mathbf{I}_n, \sigma^2) \Leftrightarrow \mathbf{y} \sim \mathcal{N}_n(\boldsymbol{\mu}, \mathbf{I}_n\sigma^2)$ and $s = \mathbf{y}'\mathbf{A}\mathbf{y} \sim \mathcal{W}_1(\nu, \sigma^2, \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu})$. Equivalently $s/\sigma^2 \sim \chi^2(\nu, \boldsymbol{\mu}'\mathbf{A}\boldsymbol{\mu}/\sigma^2)$.

19.2 DESIGN MATRIX PROPERTIES

Section 1.15 includes many useful results on functions of a design matrix \mathbf{X} of dimension $N \times q$, with $\text{rank}(\mathbf{X}) = r \leq q < N$. With \mathbf{s}_1 an $r \times 1$ vector of strictly positive values, \mathbf{L}_1 an $N \times r$ and columnwise orthonormal matrix $\mathbf{L}'_1\mathbf{L}_1 = \mathbf{I}_r$, and \mathbf{R}_1 a $q \times r$ and orthonormal matrix $\mathbf{R}'_1\mathbf{R}_1 = \mathbf{I}_r$, the SVD gives

$$\mathbf{X} = [\mathbf{L}_1 \quad \mathbf{L}_0] \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{R}'_1 \\ \mathbf{R}'_0 \end{bmatrix} = \mathbf{L}_1 \text{Dg}(\mathbf{s}_1) \mathbf{R}'_1. \tag{19.4}$$

Without loss of generality, the decomposition may be chosen such that

$$\mathbf{X}'\mathbf{X} = \mathbf{R}\text{Dg}(\mathbf{s}_1, \mathbf{0})^2 \mathbf{R}' = \mathbf{R}_1 \text{Dg}(\mathbf{s}_1)^2 \mathbf{R}'_1 \tag{19.5}$$

and

$$\mathbf{X}\mathbf{X}' = \mathbf{L} \begin{bmatrix} \text{Dg}(\mathbf{s}_1) & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{0} \end{bmatrix}^2 \mathbf{L}' = \mathbf{L}_1 \text{Dg}(\mathbf{s}_1)^2 \mathbf{L}'_1. \quad (19.6)$$

Even though $(\mathbf{X}'\mathbf{X}) = (\mathbf{X}'\mathbf{X})'$, a one-condition inverse $(\mathbf{X}'\mathbf{X})^-$ need not be symmetric. However, the special properties of $\mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}'$, a projection matrix, ensures uniqueness and symmetry of

$$\begin{aligned} \mathbf{H} &= \mathbf{X}(\mathbf{X}'\mathbf{X})^- \mathbf{X}' \\ &= \mathbf{X}(\mathbf{X}'\mathbf{X})^+ \mathbf{X}' \\ &= \mathbf{L}_1 \mathbf{L}'_1. \end{aligned} \quad (19.7)$$

The matrix \mathbf{H} is rank r and idempotent. In turn, $\mathbf{I}_N - \mathbf{H}$ is rank $\nu_e = N - r$ and idempotent. Furthermore

$$\begin{aligned} \mathbf{I}_N - \mathbf{H} &= \mathbf{I}_N - \mathbf{L}_1 \mathbf{L}'_1 \\ &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{N-r} \end{bmatrix} [\mathbf{L}_1 \ \mathbf{L}_0]' - [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \mathbf{I}_r & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} [\mathbf{L}_1 \ \mathbf{L}_0]' \\ &= [\mathbf{L}_1 \ \mathbf{L}_0] \begin{bmatrix} \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{N-r} \end{bmatrix} [\mathbf{L}_1 \ \mathbf{L}_0]' \\ &= \mathbf{L}_0 \mathbf{L}'_0. \end{aligned} \quad (19.8)$$

Lemma 19.1 In the multivariate GLM, whenever \mathbf{X} has full rank and $\mathbf{1}_N$ is the first column, the matrix $(\mathbf{X}'\mathbf{X})^{-1}$ has a simple form in terms of the first two moments of the remaining $q - 1$ predictors, \mathbf{X}_2 . With $\bar{\mathbf{x}} = \mathbf{X}'_2 \mathbf{1}_N / N$ and $\mathbf{S}_X = (\mathbf{X}'_2 \mathbf{X}_2 - N \bar{\mathbf{x}} \bar{\mathbf{x}}') / N$,

$$(\mathbf{X}'\mathbf{X})^{-1} = N^{-1} \begin{bmatrix} 1 + \bar{\mathbf{x}}' \mathbf{S}_X^{-1} \bar{\mathbf{x}} & \bar{\mathbf{x}}' \mathbf{S}_X^{-1} \\ \mathbf{S}_X^{-1} \bar{\mathbf{x}} & \mathbf{S}_X^{-1} \end{bmatrix}. \quad (19.9)$$

Proof. Without loss of generality, the intercept may be assumed to be the first column, which allows writing

$$\begin{aligned} \mathbf{X}'\mathbf{X} &= [\mathbf{1}_N \ \mathbf{X}_2]' [\mathbf{1}_N \ \mathbf{X}_2] \\ &= \begin{bmatrix} N & N \bar{\mathbf{x}}' \\ N \bar{\mathbf{x}} & \mathbf{X}'_2 \mathbf{X}_2 \end{bmatrix} \\ &= \begin{bmatrix} a_{11} & \mathbf{a}'_{21} \\ \mathbf{a}_{21} & \mathbf{A}_{22} \end{bmatrix}. \end{aligned} \quad (19.10)$$

Standard results on partitioned matrices (Theorem 1.16) gives

$$(\mathbf{X}'\mathbf{X})^{-1} = \begin{bmatrix} b_{11} & \mathbf{b}'_{21} \\ \mathbf{b}_{21} & \mathbf{B}_{22} \end{bmatrix}, \quad (19.11)$$

with

$$\begin{aligned}
 \mathbf{B}_{22} &= (\mathbf{A}_{22} - \mathbf{a}_{21}\mathbf{a}_{11}^{-1}\mathbf{a}'_{21})^{-1} \\
 &= (\mathbf{X}'_2\mathbf{X}_2 - N\bar{\mathbf{x}}N^{-1}N\bar{\mathbf{x}}')^{-1} \\
 &= (\mathbf{X}'_2\mathbf{X}_2 - N\bar{\mathbf{x}}\bar{\mathbf{x}}')^{-1} \\
 &= N^{-1}\mathbf{S}_X^{-1},
 \end{aligned} \tag{19.12}$$

$$\begin{aligned}
 b_{11} &= a_{11}^{-1} + a_{11}^{-1}\mathbf{a}'_{21}\mathbf{B}_{22}\mathbf{a}_{21}a_{11}^{-1} \\
 &= N^{-1} + N^{-1}N\bar{\mathbf{x}}'N^{-1}\mathbf{S}_X^{-1}N\bar{\mathbf{x}}N^{-1} \\
 &= N^{-1}(1 + \bar{\mathbf{x}}'\mathbf{S}_X^{-1}\bar{\mathbf{x}}),
 \end{aligned} \tag{19.13}$$

and

$$\begin{aligned}
 b_{12} &= -a_{11}^{-1}\mathbf{a}'_{21}\mathbf{B}_{22} \\
 &= -N^{-1}N\bar{\mathbf{x}}'N^{-1}\mathbf{S}_X^{-1} \\
 &= -N^{-1}\bar{\mathbf{x}}'\mathbf{S}_X^{-1}.
 \end{aligned} \tag{19.14}$$

□

19.3 MODEL COMPONENTS

For a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors and \mathbf{X} fixed at least conditionally,

$$\mathbf{E} \sim \mathcal{N}_{N,p}(\mathbf{0}, \mathbf{I}, \Sigma) \tag{19.15}$$

and

$$\mathbf{Y} \sim \mathcal{N}_{N,p}(\mathbf{X}\mathbf{B}, \mathbf{I}, \Sigma). \tag{19.16}$$

In the special case of the univariate model $p = 1$ and

$$\mathbf{e} \sim \mathcal{N}_{N,1}(\mathbf{0}, \mathbf{I}, \sigma^2) \sim \mathcal{N}_N(\mathbf{0}, \mathbf{I}\sigma^2) \tag{19.17}$$

and

$$\mathbf{y} \sim \mathcal{N}_{N,1}(\mathbf{X}\mathbf{B}, \mathbf{I}, \sigma^2) \sim \mathcal{N}_N(\mathbf{X}\mathbf{B}, \mathbf{I}\sigma^2). \tag{19.18}$$

19.4 PRIMARY PARAMETER AND RELATED ESTIMATORS

With a LTFR design matrix

$$\tilde{\mathbf{B}} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{Y}, \tag{19.19}$$

which reduces to $\tilde{\beta} = (\mathbf{X}'\mathbf{X})^{-}\mathbf{X}'\mathbf{y}$ if $p = 1$ (the univariate model). Choosing a particular $(\mathbf{X}'\mathbf{X})^{-}$ which meets the following condition is equivalent to choosing a two-condition generalized inverse:

$$[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{I}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]' = (\mathbf{X}'\mathbf{X})^{-1}. \quad (19.20)$$

With the condition,

$$\tilde{\mathbf{B}} \sim (\mathcal{S})\mathcal{N}_{q,p}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{B}, (\mathbf{X}'\mathbf{X})^{-1}, \Sigma] \quad (19.21)$$

$$\text{vec}(\tilde{\mathbf{B}}) \sim (\mathcal{S})\mathcal{N}_{qp}\{\text{vec}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{B}], \Sigma \otimes (\mathbf{X}'\mathbf{X})^{-1}\}. \quad (19.22)$$

The univariate model has $p = 1$ and

$$\tilde{\beta} \sim (\mathcal{S})\mathcal{N}_{q,1}[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta, (\mathbf{X}'\mathbf{X})^{-1}, \sigma^2] \quad (19.23)$$

$$\sim (\mathcal{S})\mathcal{N}_q[(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\beta, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]. \quad (19.24)$$

The relationships $\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{B} = \mathbf{X}\mathbf{B}$ and $\mathbf{H}\mathbf{H}\mathbf{H} = \mathbf{H}$ give

$$\begin{aligned} \hat{\mathbf{Y}} &= \mathbf{X}\tilde{\mathbf{B}} \\ &= \mathbf{H}\mathbf{Y} \\ &\sim \mathcal{SN}_{N,p}(\mathbf{X}\mathbf{B}, \mathbf{H}, \Sigma) \end{aligned} \quad (19.25)$$

$$\begin{aligned} \hat{\mathbf{E}} &= \mathbf{Y} - \hat{\mathbf{Y}} \\ &= (\mathbf{I}_N - \mathbf{H})\mathbf{Y} \\ &\sim \mathcal{SN}_{N,p}\{[\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{X}\mathbf{B}, (\mathbf{I}_N - \mathbf{H})\mathbf{I}(\mathbf{I}_N - \mathbf{H}), \Sigma\} \\ &\sim \mathcal{SN}_{N,p}[\mathbf{0}, (\mathbf{I}_N - \mathbf{H}), \Sigma]. \end{aligned} \quad (19.26)$$

The univariate case has $p = 1$,

$$\begin{aligned} \hat{y} &= \mathbf{X}\tilde{\beta} \\ &= \mathbf{H}\mathbf{y} \\ &\sim \mathcal{SN}_{N,1}(\mathbf{X}\beta, \mathbf{H}, \sigma^2) \\ &\sim \mathcal{SN}_N(\mathbf{X}\beta, \mathbf{H}\sigma^2), \end{aligned} \quad (19.27)$$

and

$$\begin{aligned} \hat{\mathbf{E}} &= \mathbf{y} - \hat{y} \\ &= (\mathbf{I}_N - \mathbf{H})\mathbf{y} \\ &\sim \mathcal{SN}_{N,1}[\mathbf{0}, (\mathbf{I}_N - \mathbf{H}), \sigma^2] \\ &\sim \mathcal{SN}_N[\mathbf{0}, (\mathbf{I}_N - \mathbf{H})\sigma^2]. \end{aligned} \quad (19.28)$$

With $\nu_e = N - r$ and $r = \text{rank}(\mathbf{X})$, in general

$$\begin{aligned} \nu_e \hat{\Sigma} &= \hat{\mathbf{E}}'\hat{\mathbf{E}} \\ &= \mathbf{Y}'(\mathbf{I}_N - \mathbf{H})(\mathbf{I}_N - \mathbf{H})\mathbf{Y} \\ &= \mathbf{Y}'(\mathbf{I}_N - \mathbf{H})\mathbf{Y} \\ &\sim \mathcal{W}_p(\nu_e, \Sigma), \end{aligned} \quad (19.29)$$

and for $p = 1$

$$\begin{aligned} \nu_e \hat{\sigma}^2 &= \hat{\mathbf{e}}' \hat{\mathbf{e}} & (19.30) \\ &= SSE \\ &= \mathbf{y}'(\mathbf{I}_N - \mathbf{H})\mathbf{y} \\ &\sim \mathcal{W}_1(\nu_e, \sigma^2). \end{aligned}$$

Equivalently,

$$\nu_e \hat{\sigma}^2 / \sigma^2 \sim \chi^2(\nu_e). \tag{19.31}$$

19.5 SECONDARY PARAMETER ESTIMATION

As always, $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'$. Having an estimable Θ requires $\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X} = \mathbf{C}$. Hence for estimable Θ

$$\begin{aligned} \hat{\Theta} &= \mathbf{C}\tilde{\mathbf{B}}\mathbf{U} & (19.32) \\ &= [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{Y}\mathbf{U}, \end{aligned}$$

$$\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\mathbf{B}\mathbf{U} = \mathbf{C}\mathbf{B}\mathbf{U}, \tag{19.33}$$

and

$$\begin{aligned} [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'][\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']' &= \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{C}' \\ &= \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\{\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{X}\}' \\ &= \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'. \end{aligned} \tag{19.34}$$

Using the results just described gives

$$\hat{\Theta} - \Theta_0 \sim \mathcal{N}_{a,b}[\mathbf{C}\mathbf{B}\mathbf{U} - \Theta_0, \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}', \Sigma] \tag{19.35}$$

$$\text{vec}(\hat{\Theta} - \Theta_0) \sim \mathcal{N}_{ab}[\text{vec}(\mathbf{C}\mathbf{B}\mathbf{U} - \Theta_0), \Sigma_* \otimes \mathbf{M}] \tag{19.36}$$

and

$$\nu_e \hat{\Sigma}_* = (N - r)\mathbf{U}'\hat{\Sigma}\mathbf{U} \sim \mathcal{W}_b(\nu_e, \mathbf{U}'\Sigma\mathbf{U}, \mathbf{0}). \tag{19.37}$$

The univariate special case has $p = 1$ and

$$\begin{aligned} \hat{\theta} &= \mathbf{C}\tilde{\beta} & (19.38) \\ &= [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}']\mathbf{y}. \end{aligned}$$

In turn,

$$\hat{\theta} - \theta_0 \sim \mathcal{N}_{a,1}[\mathbf{C}\beta - \theta_0, \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}', \sigma^2] \tag{19.39}$$

$$\sim \mathcal{N}_a(\mathbf{C}\beta - \theta_0, \sigma^2\mathbf{M}) \tag{19.40}$$

and $\nu_e \hat{\Sigma}_*$ reduces to $\nu_e \hat{\sigma}^2 = \hat{\mathbf{e}}' \hat{\mathbf{e}}$ with $\nu_e \hat{\sigma}^2 / \sigma^2 \sim \chi^2(\nu_e, 0)$.

A *testable* hypothesis requires estimable Θ , full (row) rank of C , and full (column) rank of U . The requirements ensure full rank of the $a \times a$ symmetric matrix $M = C(X'X)^-C'$ as well as nonsingular $\Sigma_* = U'\Sigma U$. In turn,

$$\begin{aligned} M &= C(X'X)^-C' \\ &= V_M \text{Dg}(\lambda_M) V_M' \\ &= [V_M \text{Dg}(\lambda_M)^{1/2}] [\text{Dg}(\lambda_M)^{1/2} V_M'] \\ &= F_M F_M' \end{aligned} \tag{19.41}$$

and

$$\begin{aligned} \widehat{\Delta} &= (\widehat{\Theta} - \Theta_0)' M^{-1} (\widehat{\Theta} - \Theta_0) \\ &= (\widehat{\Theta} - \Theta_0)' F_M^{-t} F_M^{-1} (\widehat{\Theta} - \Theta_0) \\ &= [F_M^{-1} (\widehat{\Theta} - \Theta_0)]' [F_M^{-1} (\widehat{\Theta} - \Theta_0)]. \end{aligned} \tag{19.42}$$

Furthermore

$$\begin{aligned} F_M^{-1} (\widehat{\Theta} - \Theta_0) &\sim \mathcal{N}_{a,b} [F_M^{-1} (\Theta - \Theta_0), F_M^{-1} M F_M^{-t}, \Sigma_*] \\ &\sim \mathcal{N}_{a,b} [F_M^{-1} (\Theta - \Theta_0), F_M^{-1} F_M F_M' F_M^{-t}, \Sigma_*] \\ &\sim \mathcal{N}_{a,b} [F_M^{-1} (\Theta - \Theta_0), I_a, \Sigma_*]. \end{aligned} \tag{19.43}$$

The row covariance structure being I_a allows concluding

$$\begin{aligned} \widehat{\Delta} &= [F_M^{-1} (\widehat{\Theta} - \Theta_0)]' [F_M^{-1} (\widehat{\Theta} - \Theta_0)] \\ &\sim \mathcal{W}_b \{ a, \Sigma_*, [F_M^{-1} (\Theta - \Theta_0)]' F_M^{-1} (\Theta - \Theta_0) \} \\ &\sim \mathcal{W}_b [a, \Sigma_*, (\Theta - \Theta_0)' M^{-1} (\Theta - \Theta_0)]. \end{aligned} \tag{19.44}$$

A more traditional derivation of the last result follows from considering

$$Y_U = YU - XC'(CC')^{-1}\Theta_0. \tag{19.45}$$

Having $Y \sim \mathcal{N}_{N,p}(XB, I_N, \Sigma)$ implies $Y_U \sim \mathcal{N}_{N,b}[E(Y_U), I_N, \Sigma_*]$ with

$$E(Y_U) = XBU - XC'(CC')^{-1}\Theta_0 \tag{19.46}$$

$$\Sigma_* = U'\Sigma U. \tag{19.47}$$

If $T = C(X'X)^-X'$, then $\widehat{\Theta} = TYU$. A testable Θ ensures $TXC'(CC')^{-1} = I_a$. Therefore

$$\begin{aligned} \widehat{\Theta} - \Theta_0 &= TYU - T[XC'(CC')^{-1}]\Theta_0 \\ &= TYU. \end{aligned} \tag{19.48}$$

Furthermore

$$\begin{aligned} \widehat{\Delta} &= (\widehat{\Theta} - \Theta_0)' M^{-1} (\widehat{\Theta} - \Theta_0) \\ &= Y_U' T' M^{-1} T Y_U. \end{aligned} \tag{19.49}$$

If $\mathbf{A} = \mathbf{T}'\mathbf{M}^{-1}\mathbf{T}$ then $\widehat{\Delta} = \mathbf{Y}'_U \mathbf{A} \mathbf{Y}_U$ and

$$\begin{aligned} \mathbf{A}^2 &= (\mathbf{T}'\mathbf{M}^{-1}\mathbf{T})(\mathbf{T}'\mathbf{M}^{-1}\mathbf{T}) \\ &= \mathbf{T}'\mathbf{M}^{-1}\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\{\mathbf{X}[(\mathbf{X}'\mathbf{X})^{-1}]'\mathbf{C}'\}\mathbf{M}^{-1}\mathbf{T} \\ &= \mathbf{T}'\mathbf{M}^{-1}\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\{\mathbf{C}'\}\mathbf{M}^{-1}\mathbf{T} = \mathbf{A}. \end{aligned} \tag{19.50}$$

Furthermore, with $\boldsymbol{\Theta}_{c_0} = \mathbf{C}'(\mathbf{C}\mathbf{C}')^{-1}\boldsymbol{\Theta}_0$,

$$\begin{aligned} [\mathbf{E}(\mathbf{Y}_U)]'\mathbf{A}\mathbf{E}(\mathbf{Y}_U) &= [\mathbf{X}\mathbf{B}\mathbf{U} - \mathbf{X}\boldsymbol{\Theta}_{c_0}]'\mathbf{T}'\mathbf{M}^{-1}\mathbf{T}[\mathbf{X}\mathbf{B}\mathbf{U} - \mathbf{X}\boldsymbol{\Theta}_{c_0}] \\ &= [\mathbf{T}\mathbf{X}\mathbf{B}\mathbf{U} - \mathbf{T}\mathbf{X}\boldsymbol{\Theta}_{c_0}]'\mathbf{M}^{-1}[\mathbf{T}\mathbf{X}\mathbf{B}\mathbf{U} - \mathbf{T}\mathbf{X}\boldsymbol{\Theta}_{c_0}] \\ &= (\boldsymbol{\Theta} - \boldsymbol{\Theta}_0)'\mathbf{M}^{-1}(\boldsymbol{\Theta} - \boldsymbol{\Theta}_0). \end{aligned} \tag{19.51}$$

The idempotency of \mathbf{A} combines with $\mathbf{Y}_U \sim \mathcal{N}_{N,b}[\mathbf{E}(\mathbf{Y}_U), \mathbf{I}_N, \boldsymbol{\Sigma}_*]$ to give $\widehat{\Delta} \sim \mathcal{W}_b[a, \boldsymbol{\Sigma}_*, (\boldsymbol{\Theta} - \boldsymbol{\Theta}_0)'\mathbf{M}^{-1}(\boldsymbol{\Theta} - \boldsymbol{\Theta}_0)]$ (Theorem 10.8).

19.6 ADDED-LAST AND ADDED-IN-ORDER TESTS

Muller and Fetterman (2002) gave detailed discussions of the interpretations and uses of added-last and added-in-order tests in univariate models. Here we prove some important properties.

Lemma 19.2 In a $\text{GLM}_{N,p,q}\text{FR}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors, the added-last SS for all predictors are independent if and only if $\mathbf{X}'\mathbf{X}$ is a diagonal matrix.

Proof. Testing $\boldsymbol{\beta} = \mathbf{0}$ uses $\mathbf{C} = \mathbf{I}_q$, while testing a particular slope uses $\mathbf{C}_j = \text{row}_j(\mathbf{C})$. With $\langle (\mathbf{X}'\mathbf{X})^{-1} \rangle_{jk}$ the j, k element of $(\mathbf{X}'\mathbf{X})^{-1}$,

$$m_j = \mathbf{C}_j(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'_j = [0 \ 0 \ \cdots \ 1 \ \cdots \ 0](\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'_j = \langle (\mathbf{X}'\mathbf{X})^{-1} \rangle_{jj} \tag{19.52}$$

and

$$SSH_j = \widehat{\boldsymbol{\theta}}'_j [\mathbf{C}_j(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'_j]^{-1} \widehat{\boldsymbol{\theta}}_j = \widehat{\boldsymbol{\theta}}^2_j m_j^{-1}. \tag{19.53}$$

Furthermore $\widehat{\boldsymbol{\theta}}_j = \mathbf{C}_j \widehat{\boldsymbol{\beta}}$ and

$$q_j = \widehat{\boldsymbol{\theta}}^2_j m_j^{-1} \widehat{\boldsymbol{\theta}}_j = \widehat{\boldsymbol{\beta}}' \mathbf{C}'_j m_j^{-1} \mathbf{C}_j \widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}' \mathbf{A}_j \widehat{\boldsymbol{\beta}}. \tag{19.54}$$

Here $\widehat{\boldsymbol{\beta}} \sim \mathcal{N}_q[\boldsymbol{\beta}, \sigma^2(\mathbf{X}'\mathbf{X})^{-1}]$. With $j \neq k$, Theorem 8.13 gives $q_j \perp\!\!\!\perp q_k \Leftrightarrow \mathbf{A}_j[\sigma^2(\mathbf{X}'\mathbf{X})^{-1}]\mathbf{A}_k = \mathbf{0}$. The scalar $\mathbf{C}_j(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}'_k = \langle (\mathbf{X}'\mathbf{X})^{-1} \rangle_{jk}$ gives

$$\begin{aligned} \mathbf{A}_j \sigma^2(\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}_k &= \sigma^2 \mathbf{C}'_j m_j^{-1} [\mathbf{C}_j(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_k] m_k^{-1} \mathbf{C}_k \\ &= (\sigma^2 m_j^{-1} m_k^{-1}) (\mathbf{C}'_j \mathbf{C}_k) [\mathbf{C}_j(\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_k] \\ &= (\sigma^2 m_j^{-1} m_k^{-1}) (\mathbf{C}'_j \mathbf{C}_k) \langle (\mathbf{X}'\mathbf{X})^{-1} \rangle_{jk}. \end{aligned} \tag{19.55}$$

Given the assumptions, the scalar $(\sigma^2 m_j^{-1} m_k^{-1})$ is never zero. Furthermore $(\mathbf{C}'_j \mathbf{C}_k)$ is a $q \times q$ matrix with a 1 in location j, k and zeros elsewhere. Hence the last expression is zero if and only if $\langle (\mathbf{X}'\mathbf{X})^{-1} \rangle_{jk}$ is zero. The proof is completed by noting $(\mathbf{X}'\mathbf{X})^{-1}$ is diagonal if and only if $(\mathbf{X}'\mathbf{X})$ is diagonal. \square

Lemma 19.3 In a $\text{GLM}_{N,p,q}\text{FR}(\mathbf{Y}_i; \mathbf{X}_i \mathbf{B}, \boldsymbol{\Sigma})$ with Gaussian errors, if \mathbf{X} includes an intercept, then added-last SS for predictor variables other than the intercept are mutually independent if and only if they are mutually uncorrelated.

Proof. Without loss of generality, we may assume the intercept is the leftmost column in \mathbf{X} . Testing all slopes equal to zero uses $\mathbf{C} = [\mathbf{0} \ \mathbf{I}_{q-1}]$, while testing a particular slope corresponds to considering $\mathbf{C}_j = \text{row}_j(\mathbf{C})$. Lemma 19.1 gives

$$(\mathbf{X}'\mathbf{X})^{-1} = \frac{1}{N} \begin{bmatrix} 1 + \bar{\mathbf{x}}' \mathbf{C}_X^{-1} \bar{\mathbf{x}} & \bar{\mathbf{x}}' \mathbf{C}_X^{-1} \\ \mathbf{C}_X^{-1} \bar{\mathbf{x}} & \mathbf{C}_X^{-1} \end{bmatrix}. \quad (19.56)$$

With $\langle \mathbf{C}_X^{-1} \rangle_{jk}$ indicating element j, k of \mathbf{C}_X^{-1} ,

$$m_j = \mathbf{C}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_j = [0 \ 0 \ \cdots \ 1 \ \cdots \ 0] (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}_j = N^{-1} \langle \mathbf{C}_X^{-1} \rangle_{jj} \quad (19.57)$$

and

$$SSH_j = \hat{\boldsymbol{\theta}}'_j \left[\mathbf{C}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_j \right]^{-1} \hat{\boldsymbol{\theta}}_j = \hat{\boldsymbol{\theta}}_j^2 m_j^{-1}. \quad (19.58)$$

Furthermore $\hat{\boldsymbol{\theta}}_j = \mathbf{C}_j \hat{\boldsymbol{\beta}}$,

$$q_j = \hat{\boldsymbol{\theta}}'_j m_j^{-1} \hat{\boldsymbol{\theta}}_j = \hat{\boldsymbol{\beta}}' \mathbf{C}'_j m_j^{-1} \mathbf{C}_j \hat{\boldsymbol{\beta}} = \hat{\boldsymbol{\beta}}' \mathbf{A}_j \hat{\boldsymbol{\beta}}. \quad (19.59)$$

Here $\hat{\boldsymbol{\beta}} \sim N_q[\boldsymbol{\beta}, \sigma^2 (\mathbf{X}'\mathbf{X})^{-1}]$. With $j \neq k$, Theorem 8.13 gives $q_j \perp\!\!\!\perp q_k \Leftrightarrow \mathbf{A}_j [\sigma^2 (\mathbf{X}'\mathbf{X})^{-1}] \mathbf{A}_k = \mathbf{0}$. The scalar $\mathbf{C}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_k = \langle \mathbf{C}_X^{-1} \rangle_{jk}$ gives

$$\begin{aligned} \mathbf{A}_j \sigma^2 (\mathbf{X}'\mathbf{X})^{-1} \mathbf{A}_k &= \sigma^2 \mathbf{C}'_j m_j^{-1} \left[\mathbf{C}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_k \right] m_k^{-1} \mathbf{C}_k \\ &= (\sigma^2 m_j^{-1} m_k^{-1}) (\mathbf{C}'_j \mathbf{C}_k) \left[\mathbf{C}_j (\mathbf{X}'\mathbf{X})^{-1} \mathbf{C}'_k \right] \\ &= (\sigma^2 m_j^{-1} m_k^{-1}) (\mathbf{C}'_j \mathbf{C}_k) \langle \mathbf{C}_X^{-1} \rangle_{jk}. \end{aligned} \quad (19.60)$$

Given the assumptions, the scalar $(\sigma^2 m_j^{-1} m_k^{-1})$ is never zero. Furthermore $(\mathbf{C}'_j \mathbf{C}_k)$ is $q \times q$ with a 1 in location j, k and zeros elsewhere. Hence the last expression is zero if and only if $\langle \mathbf{C}_X^{-1} \rangle_{jk}$ is zero. The proof is completed by noting \mathbf{C}_X^{-1} is diagonal if and only if \mathbf{C}_X is diagonal if and only if the corresponding correlation matrix is diagonal. \square

Lemma 19.4 Added-in-order SS are always independent in a $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ with Gaussian errors.

Proof. In the following, \mathbf{X}_j indicates the first j columns of \mathbf{X} , and $\mathbf{X}_{i,j} = \text{row}_i(\mathbf{X}_j)$. In turn, model j is $\text{GLM}_{N,q}\text{FR}(y_i; \mathbf{X}_{i,j}\boldsymbol{\beta}_j, \sigma_j^2)$ with Gaussian errors. Furthermore $\mathbf{y} = \mathbf{X}_j\boldsymbol{\beta}_j + \mathbf{e}_j$, with full-rank \mathbf{X}_j of dimension $N \times j$. Also

$$\begin{aligned} \mathbf{X}_1\boldsymbol{\beta}_1 &= [\mathbf{x}_1][\beta_1] \\ &\vdots \\ \mathbf{X}_j\boldsymbol{\beta}_j &= [\mathbf{x}_1 \ \mathbf{x}_2 \ \cdots \ \mathbf{x}_j] \begin{bmatrix} \beta_1 \\ \vdots \\ \beta_j \end{bmatrix}. \end{aligned} \tag{19.61}$$

The matrix $\mathbf{H}_j = \mathbf{X}_j(\mathbf{X}'_j\mathbf{X}_j)^{-1}\mathbf{X}'_j = \mathbf{H}_j^2$ is $N \times N$, symmetric, idempotent, and $\text{rank}(\mathbf{H}_j) = \text{rank}(\mathbf{X}_j) = j$. Similarly, $(\mathbf{I}_N - \mathbf{H}_j)$ is symmetric, $N \times N$, idempotent, of rank $N - j$, and $\mathbf{H}_j(\mathbf{I}_N - \mathbf{H}_j) = (\mathbf{I}_N - \mathbf{H}_j)\mathbf{H}_j = \mathbf{0}$. Also

$$\begin{aligned} \mathbf{y}'\mathbf{y} &= \mathbf{y}'\mathbf{I}_N\mathbf{y} = \mathbf{y}'\mathbf{H}_j\mathbf{y} + \mathbf{y}'(\mathbf{I}_N - \mathbf{H}_j)\mathbf{y} \\ \text{SST}_j &= \text{SSH}_j + \text{SSE}_j. \end{aligned} \tag{19.62}$$

The ANOVA theorem may be applied by defining $\mathbf{y}_{1,j} = \mathbf{y}\sigma_j^{-1}$ and writing

$$\mathbf{A} = \mathbf{A}_1 + \mathbf{A}_2 \tag{19.63}$$

$$\mathbf{I}_N = \mathbf{H}_j + (\mathbf{I}_N - \mathbf{H}_j), \tag{19.64}$$

with a corresponding decomposition of ranks, namely $N = j + (N - j)$.

The test of adding variables $\{j + 1, j + 2, \dots, j + k\}$ in order compares model j to model $j + k$ (models in the added-in-order pool) and uses the added-in-order sums of squares, $\text{SS}_{j+k} = \text{SSE}_j - \text{SSE}_{j+k} = \text{SSH}_{j+k} - \text{SSH}_j$. In particular,

$$\text{SS}_{j+k} = \mathbf{y}'\mathbf{H}_{j+k}\mathbf{y} - \mathbf{y}'\mathbf{H}_j\mathbf{y} = \mathbf{y}'(\mathbf{H}_{j+k} - \mathbf{H}_j)\mathbf{y} = \mathbf{y}'\mathbf{A}_{j+k}\mathbf{y}. \tag{19.65}$$

The design matrix for model $j + k$ may be partitioned as $\mathbf{X}_{j+k} = [\mathbf{X}_j \ \mathbf{X}_{k(-j)}]$, with j columns in \mathbf{X}_j and $k - j$ columns in $\mathbf{X}_{k(-j)}$, which contains variables $\{j + 1, j + 2, \dots, j + k\}$. In turn,

$$\mathbf{X}'_{j+k}\mathbf{X}_{j+k} = \begin{bmatrix} \mathbf{X}'_j\mathbf{X}_j & \mathbf{X}'_j\mathbf{X}_{k(-j)} \\ \mathbf{X}'_{k(-j)}\mathbf{X}_j & \mathbf{X}'_{k(-j)}\mathbf{X}_{k(-j)} \end{bmatrix} \tag{19.66}$$

has a similarly partitioned inverse,

$$(\mathbf{X}'_{j+k}\mathbf{X}_{j+k})^{-1} = \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}'_{12} & \mathbf{B}_{22} \end{bmatrix}, \tag{19.67}$$

with

$$\mathbf{B}_{11} = (\mathbf{X}'_j\mathbf{X}_j)^{-1} + (\mathbf{X}'_j\mathbf{X}_j)^{-1}\mathbf{X}'_j\mathbf{X}_{k(-j)}\mathbf{B}_{22}\mathbf{X}'_{k(-j)}\mathbf{X}_j(\mathbf{X}'_j\mathbf{X}_j)^{-1} \tag{19.68}$$

$$\mathbf{B}_{12} = -(\mathbf{X}'_j\mathbf{X}_j)^{-1}\mathbf{X}'_j\mathbf{X}_{k(-j)}\mathbf{B}_{22}. \tag{19.69}$$

Using the partitioned matrix inverse in the hat matrix for the larger model gives

$$\begin{aligned}
 H_{j+k} &= \mathbf{X}_{j+k} (\mathbf{X}'_{j+k} \mathbf{X}_{j+k})^{-1} \mathbf{X}'_{j+k} \\
 &= [\mathbf{X}_j \mathbf{X}_{k(-j)}] \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}'_{12} & \mathbf{B}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{X}'_j \\ \mathbf{X}'_{k(-j)} \end{bmatrix} \\
 &= [\mathbf{X}_j \mathbf{X}_{k(-j)}] \begin{bmatrix} \mathbf{B}_{11} & \mathbf{B}_{12} \\ \mathbf{B}'_{12} & \mathbf{B}_{22} \end{bmatrix} \begin{bmatrix} \mathbf{X}'_j \\ \mathbf{X}'_{k(-j)} \end{bmatrix} \\
 &= [\mathbf{X}_j \mathbf{B}_{11} + \mathbf{X}_{k(-j)} \mathbf{B}'_{12} \quad \mathbf{X}_j \mathbf{B}_{12} + \mathbf{X}_{k(-j)} \mathbf{B}_{22}] \begin{bmatrix} \mathbf{X}'_j \\ \mathbf{X}'_{k(-j)} \end{bmatrix} \\
 &= \mathbf{X}_j \mathbf{B}_{11} \mathbf{X}'_j + \mathbf{X}_{k(-j)} \mathbf{B}'_{12} \mathbf{X}_j + \mathbf{X}_j \mathbf{B}_{12} \mathbf{X}'_{k(-j)} + \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}. \quad (19.70)
 \end{aligned}$$

The first term in the last line of the preceding equation may be expressed as

$$\begin{aligned}
 \mathbf{X}_j \mathbf{B}_{11} \mathbf{X}'_j &= \mathbf{X}_j (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j + \\
 &\quad \mathbf{X}_j (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{X}_j (\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j \\
 &= \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j, \quad (19.71)
 \end{aligned}$$

while the third term may be expressed as

$$\begin{aligned}
 \mathbf{X}_j \mathbf{B}_{12} \mathbf{X}'_{k(-j)} &= \mathbf{X}_j [-(\mathbf{X}'_j \mathbf{X}_j)^{-1} \mathbf{X}'_j \mathbf{X}_{k(-j)} \mathbf{B}_{22}] \mathbf{X}'_{k(-j)} \\
 &= -\mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}. \quad (19.72)
 \end{aligned}$$

Substituting the alternate forms back into the expression for \mathbf{H}_{j+k} gives

$$\begin{aligned}
 H_{j+k} &= \mathbf{X}_j \mathbf{B}_{11} \mathbf{X}'_j + \mathbf{X}_{k(-j)} \mathbf{B}'_{12} \mathbf{X}'_j + \mathbf{X}_j \mathbf{B}_{12} \mathbf{X}'_{k(-j)} + \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \\
 &= \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j - \\
 &\quad (\mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}) + \\
 &\quad \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}. \quad (19.73)
 \end{aligned}$$

In turn,

$$\begin{aligned}
 \mathbf{H}_j \mathbf{H}_{j+k} &= \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j - \\
 &\quad (\mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}) + \\
 &\quad \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \\
 &= \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j - \\
 &\quad (\mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} \mathbf{H}_j + \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)}) + \\
 &\quad \mathbf{H}_j \mathbf{X}_{k(-j)} \mathbf{B}_{22} \mathbf{X}'_{k(-j)} = \mathbf{H}_j. \quad (19.74)
 \end{aligned}$$

Also $\mathbf{H}_{j+k} = \mathbf{H}_j + (\mathbf{H}_{j+k} - \mathbf{H}_j) = \mathbf{H}_j + \mathbf{A}_{j+k}$, $\mathbf{H}_j \mathbf{A}_{j+k} = \mathbf{H}_j (\mathbf{H}_{j+k} - \mathbf{H}_j) = \mathbf{0}$, and $\text{rank}(\mathbf{H}_{j+k}) = \text{rank}(\mathbf{H}_j) + \text{rank}(\mathbf{A}_{j+k})$. Finally

$$\begin{aligned}
 \mathbf{A}_{j+k_1} \mathbf{A}_{(j+k_1)+k_2} &= (\mathbf{H}_{j+k_1} - \mathbf{H}_j)(\mathbf{H}_{j+k_1+k_2} - \mathbf{H}_{j+k_1}) \\
 &= \mathbf{H}_{j+k_1} \mathbf{H}_{j+k_1+k_2} - \mathbf{H}_j \mathbf{H}_{j+k_1+k_2} - \mathbf{H}_{j+k_1} \mathbf{H}_{j+k_1} + \mathbf{H}_j \mathbf{H}_{j+k_1} \\
 &= \mathbf{H}_{j+k_1} - \mathbf{H}_j - \mathbf{H}_{j+k_1} + \mathbf{H}_j \\
 &= \mathbf{0}
 \end{aligned}
 \tag{19.75}$$

and

$$\begin{aligned}
 \mathbf{A}_{j+k_1} \mathbf{A}_{(j+k_1+m)+k_2} &= (\mathbf{H}_{j+k_1} - \mathbf{H}_j)(\mathbf{H}_{j+k_1+m+k_2} - \mathbf{H}_{j+k_1+m}) \\
 &= \mathbf{H}_{j+k_1} \mathbf{H}_{j+k_1+m+k_2} - \mathbf{H}_j \mathbf{H}_{j+k_1+m+k_2} - \\
 &\quad \mathbf{H}_{j+k_1} \mathbf{H}_{j+k_1+m} + \mathbf{H}_j \mathbf{H}_{j+k_1+m} \\
 &= \mathbf{H}_{j+k_1} - \mathbf{H}_j - \mathbf{H}_{j+k_1} + \mathbf{H}_j \\
 &= \mathbf{0}.
 \end{aligned}
 \tag{19.76}$$

□

Sample Size for Univariate Linear Models

20.1 SAMPLE SIZE CONSULTING: BEFORE YOU BEGIN

A scientist describes, in three sentences, a study requiring six months and costing tens of thousands of dollars. The scientist then asks “How many subjects do I need?” Finding a good answer requires intensive and iterative collaboration between the scientist and statistician to clearly detail the (1) ethical, monetary, and time constraints, (2) scientific goals, both long term (vague) and short term (concrete), (3) study design, and (4) data analysis plan. We believe some form of sample size analysis, such as power analysis for test procedures and analysis of precision for estimators, should play a key role in the planning of most studies. Muller, Barton, and Benignus (1984), Muller and Benignus (1992), O'Brien and Muller (1993), and Catellier and Muller (2002) provided introductory and tutorial presentations.

In practice, scientists typically seek to achieve more than one goal in each study. Finding a satisfactory design and analysis plan requires eliciting the complete set of goals and the relative importance of each. Although estimation essentially always has importance, the importance of statistical hypothesis testing can range from negligible (or not applicable) to critically essential. An important question for the statistician to ask the scientist is, “Will the study be a success if none of the planned hypothesis tests turn out to be statistically significant?” Jiroutek, Muller, Kupper, and Stewart (2003) discussed the question from the perspective of choosing a criterion probability to control when choosing a sample size.

In addition to choosing a sample size, the collaborative process should include comparing a variety of designs and associated analyses. Exploring and describing the variation in performance as a function of design features, assumptions about the population, and choice of analysis inform and improve the decision process.

Tables and plots for a range of different scenarios help greatly. Figure 20.1 illustrates the power tradeoffs between sample size and mean difference for an independent groups t -test. As with all univariate linear models with Gaussian errors and fixed predictors, given α , only 1) sample size, 2) mean differences and 3) error variance affect power. Example E01 in the manual for the free software

POWERLIB203.IML described in the Appendix (Section A.2) gives the code used to produce Figure 20.1.

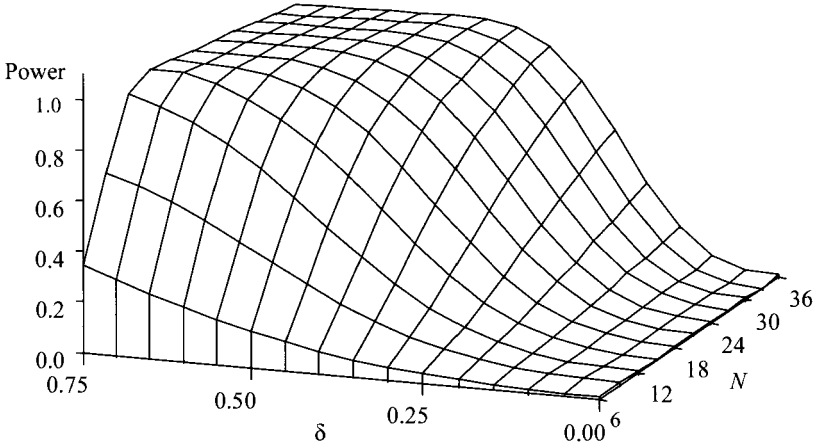


Figure 20.1 Power as a function of sample size (N) and mean difference (δ) for a balanced independent groups t test with $\sigma^2 = 0.068$ and $\alpha = 0.01$.

20.2 THE MACHINERY OF A POWER ANALYSIS

The discussion assumes Gaussian errors and fixed predictors. We restrict attention to testable hypotheses, which requires full-rank $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X})$. A univariate GLM power calculation with fixed predictors is fully specified by α , σ^2 , \mathbf{X} , $\boldsymbol{\beta}$, \mathbf{C} , and $\boldsymbol{\theta}_0$. The size of the test, α , as well as the test statistic must be chosen a priori. The dimensions of the model fix the degrees of freedom. Although specifying $\boldsymbol{\beta}$ and σ^2 suffices to complete the power analysis, specifying the (usually smaller and simpler) matrices $\boldsymbol{\theta} = \mathbf{C}\boldsymbol{\beta}$, $\boldsymbol{\theta}_0$, and σ^2 also suffices and usually proves easier. A further simplification occurs because the noncentrality $\omega = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0) / \sigma^2$ suffices. With $\omega = N\rho^2 / (1 - \rho^2)$, ρ^2 is a (generalized) squared correlation. Also $0 \leq \omega < \infty$ and $0 \leq \rho^2 \leq 1$, while $\omega \equiv 0 \Leftrightarrow \rho^2 \equiv 0 \Leftrightarrow \boldsymbol{\theta} = \boldsymbol{\theta}_0 \Leftrightarrow H_0 \text{ holds} \Leftrightarrow \text{power} = \alpha$.

Deleting any duplicate rows from the design matrix creates the essence matrix (Definition 11.5), $\text{Es}(\mathbf{X})$. It allows easily determining essential properties of a design, such as rank. Comparing essence matrices allows determining relationships between alternate parameter definitions. The concept allows conveniently separating total sample size from the coding scheme. The separation simplifies computing and interpreting power in linear models.

20.3 INDEPENDENT *t* EXAMPLE

For cell mean coding, $Es(\mathbf{X}) = \mathbf{I}_2$, $\mathbf{C} = [1 \ -1]$, and $\mathbf{U} = 1$, while for reference cell coding $Es(\mathbf{X}) = \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix}$ and $\mathbf{C} = [0 \ 1]$. Here n indicates the number of replicates (number of observations per unique row of \mathbf{X} , which corresponds to a cell for any factorial design). In turn,

$$\begin{aligned} \mathbf{y} &= (\mathbf{1}_n \otimes \mathbf{I}_2) \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \mathbf{e} \\ &= \left(\mathbf{1}_n \otimes \begin{bmatrix} 1 & 0 \\ 1 & 1 \end{bmatrix} \right) \begin{bmatrix} \mu \\ \delta \end{bmatrix} + \mathbf{e}. \end{aligned} \tag{20.1}$$

Here $[\mu \ \delta]'$ is equivalent to $\delta \cdot [0 \ 1]'$ in terms of μ for the test of interest.

Example 20.1 The following code computes power for the model just described.

```
TITLE1 "P0802.SAS--simple independent t test";
PROC IML SYMSIZE=4000 WORKSIZE=4000;
%INCLUDE "..\IML\POWERLIB.IML";
C={1 -1};    U={1};    THETA0=0;
ALPHA=.01;
SIGMA={2.1};*Variance if univariate, as here;  SIGSCAL={1};
RHOSCAL={1};
ESSENCEX=1(2); REPN=5;    or REPN={5,10};
BETA={7.0,7.0 };  BETASCAL={1}; RUN POWER;
BETA={7.0,7.3 };  BETASCAL={1}; RUN POWER;
BETA={7.0,7.15};  BETASCAL={1}; RUN POWER;
*last 3 lines together equivalent to either of next 2 lines;
*  BETA={0,1}; BETASCAL={0,.15,.30}; RUN POWER;
*  BETA={0,1}; BETASCAL=DO(0,.30,.15); RUN POWER;
```

Here $\theta = (\mu_1 - \mu_2) = \delta$. Hence one may choose $BETA=\{0,1\}$ and $BETASCAL=\delta$ or $BETASCAL=DO(\delta_{low}, \delta_{high}, \delta_{increment})$. Doing so avoids specifying the grand mean, $(\mu_1 + \mu_2)/2$, which does not affect the power of interest.

__	HOLDPOW	CASE	ALPHA	SIGSCAL	RHOSCAL	BETASCAL	TOTAL_N	POWER
	1	0.01	1	1	0	10	0.01	
	2	0.01	1	1	0	20	0.01	
	3	0.01	1	1	0.15	10	0.011	
	4	0.01	1	1	0.15	20	0.012	
	5	0.01	1	1	0.3	10	0.013	
	6	0.01	1	1	0.3	20	0.017	

What if an unbalanced design is required? Two methods are available.

Method 1. REPN=1.

Example with $N_1 = 11$ participants in one group and $N_2 = 7$ in the second group.

REPN=1;

N1=11; N2=7; ESSENCEX=BLOCK(J(N2,1,1), J(N1,1,1)); $\mathbf{X} = \begin{bmatrix} \mathbf{1}_7 & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{11} \end{bmatrix}$;

Method 2. REPN \geq 1.

Example proportional design with two controls for every treatment participant.

ESSENCEX=BLOCK(J(2,1,1), J(1,1,1)); $\mathbf{X} = \begin{bmatrix} 1 & 0 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}$;

Allows using REPN={3,6}; to consider N=9 and N=18.

20.4 PAIRED t EXAMPLE

A paired data t test can be conducted as a one-sample t test with the model

$$\mathbf{y} = \mathbf{1}_n \delta + \mathbf{e}. \quad (20.2)$$

Example 20.2

```
TITLE1 "P0804.SAS--paired t, one sample/difference version";
PROC IML SYMSIZE=4000 WORKSIZE=4000;
%INCLUDE "..\IML\POWERLIB.IML";
C={1}; U={1}; THETA0=0;
ALPHA=.01;
SIGMA={ 1.1 };*Variance of difference; SIGSCAL={1}; RHOSCAL={1};
ESSENCEX=I(1);
REPN={5,10};
BETA={1}; BETASCAL=DO(0,.3,.15);
RUN POWER;
```

20.5 THE IMPACT OF USING $\hat{\sigma}^2$ OR $\hat{\beta}$ IN POWER ANALYSIS

Data analysts often use $\hat{\sigma}^2$, an estimator of σ^2 , in a power analysis. In turn, the power value inherits the randomness of the estimator. Taylor and Muller (1995) derived exact methods to account for the randomness in the context of the $GLM_{N,q}(y_i; \mathbf{X}_i \boldsymbol{\beta}, \sigma^2)$ with fixed \mathbf{X} and Gaussian errors. Accounting for the randomness leads to creating a confidence region around the power curve, as illustrated in Figure 20.2. The particular range of values used in the figure reflect the following scientific context. In humans, substantially elevated creatinine levels in the blood typically indicate severe kidney disease. With normal values falling below 1 mg/dL, values of 2 or higher are important and a bad sign. Values in the range of 10–20 usually reflect kidneys near complete failure (leading to dialysis, kidney transplant, or death). Experience with the measure has led to the realization that $1/\text{creatinine}$ (in units of dL/mg) is approximately Gaussian. The vertical reference line at 0.5 dL/mg in Figure 20.2 indicates the change in reciprocal

creatinine deemed of clinical significance by the nephrologist (kidney specialist) who asked for guidance on sample size.

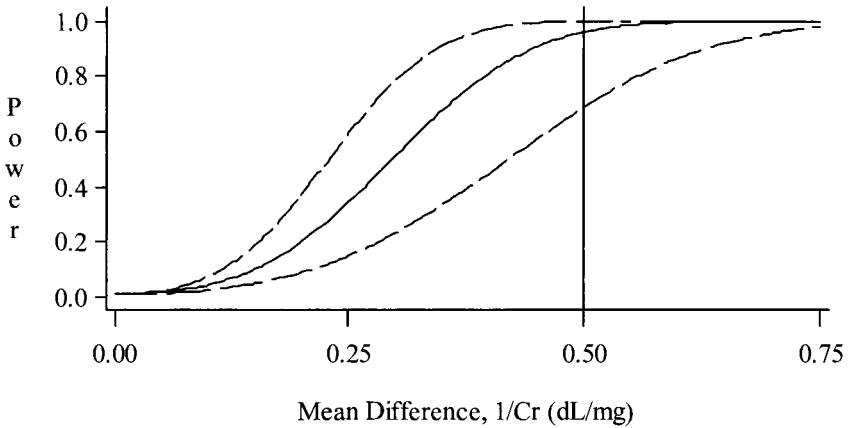


Figure 20.2 Power curve with two-sided 95% confidence region due to $\hat{\sigma}^2$ for independent groups t test power with $N = 12 + 12$.

In some cases, data analysts conduct power analysis conditional on the outcome of a study. At the end of a trial in the drug discovery process, planning for future studies, including power analysis, centers on the conditions and variables with the smallest p values. In other settings, a strong intellectual commitment to a particular hypothesis may lead a scientist to conduct a power analysis following a nonsignificant result. The first scenario implies $\hat{\sigma}^2$ arose from a density truncated on the right (excluding large values), while the second scenario implies $\hat{\sigma}^2$ arose from a density truncated on the left (excluding small values). Muller and Pasour (1997) extended the results of Taylor and Muller (1995) to allow for such truncation. The first scenario creates optimistic bias (estimated power too large), while the second scenario creates pessimistic bias (estimated power too small).

Taylor and Muller (1996) described the impact of using estimates of both σ^2 and β in estimating power of a univariate linear model. We have seen the approach used in three settings: data analysis, study planning for assessing individual response, and study planning for population response. Only the last application seems defensible.

In the context of data analysis, Lenth (2001) appropriately criticized the computation of such an estimate (“retrospective power,” \hat{P}). Taylor and Muller’s (1996) results make it clear that \hat{P} is a one-to-one function of the p value. Hence it adds no information or value to any data analysis. Finding the largest \hat{P} is equivalent to finding the smallest p value, which is often very misleading.

Lenth (2001) also appropriately criticized estimating both σ^2 and β when planning a future study (“prospective power”). Allowing an observed difference to drive the sample size makes no reference to the concept of scientific importance. In Figure 20.2, the clinically important difference of 0.5 dL/mg (1/Cr) drove the

power analysis. In the context of evaluating a drug intended to improve kidney function, with response variable $1/Cr$, it does not seem defensible to use an estimate of β to drive the power analysis. The setting involves study planning for assessing individual response.

In contrast, estimating both σ^2 and β in study planning for assessing population response may be defensible. Taylor and Muller (1996) illustrated the process in the context of U.S. EPA funded research on the effects of carbon monoxide (CO) on human perceptual-motor performance, such as driving an automobile. For a particular exposure level, the performance decrement for a single individual might be modest or even negligible. However, the same exposure level experienced by an entire population may lead to an unacceptable level of total risk. Hence EPA scientists sought to replicate the most credible, not the largest, published finding. Taylor and Muller provided a detailed description of the application.

Taylor and Muller (1996) also described the impact of conducting the power analysis condition on the outcome of the study providing the estimate. As when estimating only σ^2 , requiring the previous study to have a significant result creates optimistic bias (estimated power too large). Similarly, requiring the previous study to have a nonsignificant result creates pessimistic bias (estimated power too small).

20.6 RANDOM PREDICTORS

Although errors in measurement could introduce additional randomness in both fixed and random predictors, we assume that the scientist measures all predictors without appreciable error. Given the assumption, in practice, the distinction between random and fixed predictors does not affect the distribution theory. However, in power analysis the distinction between random and fixed predictors changes and complicates the distribution theory. Sampson (1974) detailed many of the basic issues and known results for both the univariate and multivariate model with Gaussian predictors.

Jayakar (1970), Soller and Genizi (1978), Genizi and Soller (1979), and Gatsonis and Sampson (1989) developed methods for models involving random dichotomous and Gaussian predictors in the univariate GLM. To our knowledge, power with any other predictor distribution has not been studied.

Many questions remain, especially with combinations of fixed and random predictors. Most importantly, do the simple approximations often used in practice lead to poor approximations of power?

As with any probability, power can be interpreted as the expected value of an indicator variable. Power computed with random predictors can be thought of as *expected power* due to expectation with respect to the choice of predictor values. Glueck and Muller (2003) recommended considering *quantile power*, such as median power, in lieu of expected power. A more conservative approach would use a lower quantile to reduce risk of study failure.

20.7 INTERNAL PILOT DESIGNS

Having chosen a β of scientific importance (such as a pattern of mean differences), a valid choice of σ^2 usually stands as the biggest barrier to an appropriate choice of sample size in the univariate linear model. An *internal pilot design* (Wittes and Brittain, 1990) solves the problem as follows. First, a traditional power analysis is conducted based on a planning value σ_0^2 and a target power which together imply a total sample size of n_0 . Second, the first $n_1 < n_0$ observations are collected, and $\hat{\sigma}_1^2$ is computed from residuals from the appropriate linear model. However, no interim data analysis is conducted, and the data analysts and scientists remain masked with respect to treatment assignment. Third, a new total sample size, $N_+ = n_1 + N_2$, is computed, based on $\hat{\sigma}_1^2$. Fourth, an additional N_2 observations are collected. Fifth, and finally, the analysis is conducted for the complete set of N_+ observations. The interim power analysis typically increases sample size when needed to compensate for σ_0^2 being too small and decreases sample size to compensate for σ_0^2 being too large. Hence expected sample size and power are improved.

Unfortunately, an internal pilot design can inflate test size, at least in small samples. Hence Coffey and Muller (1999, 2000a, 2000b, 2001) described many small-sample results which allow using an internal pilot with any univariate linear model with Gaussian errors and fixed predictors. The methods control test size while still providing the desired advantages in power and expected sample size.

20.8 OTHER CRITERIA FOR CHOOSING A SAMPLE SIZE

In the context of the Neyman-Pearson approach to testing a hypothesis, test size, indicated α , equals the probability of rejecting the null hypothesis given the null holds: $\alpha = \Pr\{\text{reject } H_0 | H_0 = \text{TRUE}\} = \Pr\{\text{reject } H_0 | H_A = \text{FALSE}\}$. While α equals the probability of a type I error, a false positive, $\beta = \Pr\{\text{fail to reject } H_0 | H_A = \text{TRUE}\}$ gives the probability of a type II error, a false negative. In turn, the conventional definition of power is $1 - \beta = \Pr\{\text{reject } H_0 | H_A = \text{TRUE}\}$. It is mathematically convenient to define the power function, $\Pr\{\mathbf{y} \in RR | H\}$, as a function of θ by adding a single point so that power under the null is at most α .

A goal other than rejecting a hypothesis leads to using design criteria other than power. Historically, interest in controlling confidence interval width has been the next most popular criterion. For scalar hypotheses in the general linear model, Jiroutek, Muller, Kupper and Stewart (2003) reviewed criteria in terms of three basic events. The event *width* (W) occurs when the observed confidence interval is less than a fixed constant chosen a priori. The event *validity* (V) occurs if the confidence interval contains the parameter of interest. The event *rejection* (R) occurs if the confidence interval excludes the null value of the hypothesis. For a two-sided test in the general linear model, the probability of rejection is slightly greater than power due to the slight chance that the hypothesis is rejected in the

“wrong” direction. For a t test with positive difference in population means, rejection in the wrong direction involves rejecting the null with a negative difference observed. For a one-sided test, the probability of rejection and the (unconditional) definition of power coincide.

Jiroutek et al. (2003) observed that various authors have considered controlling only power, essentially $\Pr\{R\}$, or $\Pr\{R|V\}$, or $\Pr\{W\}$, or $\Pr\{W|V\}$. They argued that scientists often desire not only a valid hypothesis test with good power but also a valid confidence interval of a reasonable size. If so, then $\Pr\{(W \cap R)|V\}$ captures the combined goals, with all previously studied criteria occurring as special cases. Jiroutek et al. described convenient single-integral formulas for computing $\Pr\{(W \cap R)|V\}$. Not surprisingly, controlling $\Pr\{W\}$ or $\Pr\{W|V\}$ when study goals include rejection can lead to collecting far too few or far too many observations. Either problem may also occur when controlling power and ignoring confidence interval goals. Free SAS/IML[®] software to compute $\Pr\{(W \cap R)|V\}$, $\Pr\{R\}$, $\Pr\{W\}$, and $\Pr\{W|V\}$ for any scalar hypothesis in a univariate or multivariate GLM with fixed predictors can be found at the website documented in the Appendix (Section A.2).

EXERCISES

20.1 Use the POWERLIB software described in the Appendix (Section A.2) to reproduce the results in Example 20.1.

20.2 A GLM _{N,q} ($y_i; \mathbf{X}_i\beta, \sigma^2$) with Gaussian errors and $N = 10$ observations has

$$\mathbf{X}\beta = \begin{bmatrix} \mathbf{1}_{10} & \mathbf{x} \end{bmatrix} \begin{bmatrix} \alpha \\ \gamma \end{bmatrix}.$$

With $\mathbf{C} = [0 \ 1]$ and $\mathbf{M} = [\mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}']$, a two-sided test of $H_0 : \gamma = 0$ uses

$$F_{\text{obs}} = (\mathbf{C}\hat{\beta})' \mathbf{M}^{-1} (\mathbf{C}\hat{\beta}) / \hat{\sigma}^2 = \frac{\hat{\gamma}^2}{\hat{\sigma}^2} \sum_{i=1}^N (x_i - \bar{x})^2 \sim F(1, N-2, \omega),$$

and noncentrality $\omega = (\mathbf{C}\beta)' \mathbf{M}^{-1} (\mathbf{C}\beta) / \sigma^2 = (\gamma^2 / \sigma^2) \sum_{i=1}^N (x_i - \bar{x})^2$.

Consider testing the null hypothesis $H_0 : \gamma = 0$ with a significance level of 0.05.

20.2.1 What is the probability of rejecting the null hypothesis if $\alpha = 1$, $\gamma = 1/4$, $\sigma^2 = 1$, and $\mathbf{x} = [0.3 \ 1.0 \ 1.5 \ 1.7 \ 1.8 \ 2.5 \ 2.6 \ 4.0 \ 7.4 \ 14.5]'$?

20.2.2 Suppose that instead of being fixed, \mathbf{x} is the realization of a vector-valued random variable. Specifically, let \mathbf{x} be a random sample of size 10 from an exponential distribution with mean $\mu_x = 4$. Use computer simulations to estimate the *unconditional* probability of rejecting the null hypothesis.

Hints. (a) $\Pr\{\text{REJECT NULL}\} = \mathbf{E}_{\mathbf{x}}(\Pr\{\text{REJECT NULL}|\mathbf{x}\}) = \mathbf{E}_{\mathbf{x}}[g(\mathbf{x})]$ (say).

(b) Let $\{x_1, \dots, x_N\}$ denote a random sample and let $y_i = g(x_i)$. The statistic $\bar{y} = \sum_{i=1}^N y_i / N$ is an unbiased estimator of the unconditional probability.

(c) The simulations can be conducted in SAS/IML or another matrix-based language. The following is a skeleton program.

```
SEED=nnnnn; *you must assign a value to start
              pseudo-random number generator;
*Ideally, choose a prime number as large as possible;
*In practice, choose 5-7 digit # not divisible by 2, 3, or 5;
*Number nmopq is divisible by 3 <=> n+m+o+p+q divisible by 3;
NREPS = 1000; * THE NUMBER OF REPLICATIONS;
power = J(NREPS,1, . ); * AN EMPTY VECTOR TO STORE RESULTS;
do i = 1 to NREPS;
  x=4*ranexp(J(10,1,SEED)); * GENERATE RANDOM SAMPLE, SIZE 10;
  NCP = ... ; * A FUNCTION OF  $\gamma$ ,  $\sigma^2$ , and x.
  CRIT = ...; * CALCULATE THE CRITICAL VALUE;
  power[i] = ...; * A FUNCTION of NCP, CRIT df1, df2;
end;
estimate = sum(power)/NREPS; * ESTIMATED UNCONDITIONAL POWER;
```

The approach suggested is an extremely inefficient programming strategy. The data should be “buffered” (written to disk after, perhaps, every 100 replications). We accept such inefficiency to help keep the exercise manageable.

20.3 A researcher plans a randomized trial to compare three drugs (H_2 receptor antagonists) used to treat stomach ulcers. An equal number of patients will be randomized to three treatment groups (Cimetidine, Ranitidine, Famotidine). The main outcome is gastric acidity measured in pH. The goal is to test whether the three drugs increase pH the same amount. There are no covariates. The investigator is not sure whether to obtain baseline measurements. She knows that baseline measurements can improve power but is not sure whether it is worth the cost. With a fixed budget, money spent on baseline measurements reduces the number of participants. She has asked help in choosing the most cost-effective and powerful design. All tests will be based on a type I error rate of $\alpha = 0.05$.

Design 1: one measurement per participant

Measure gastric acidity after a 6 month treatment regimen. Use a univariate one-way ANOVA model with 3 groups. The outcome variable is gastric acidity (call it y_1). The null hypothesis is that $E(y_1)$ is the same for all participants in all treatments. The test statistic is the ANOVA test of overall regression, which has 2 numerator degrees of freedom.

Design 2: difference scores, two measurements per participant (baseline and follow-up)

Obtain a “baseline” measurement prior to initiating therapy (call it y_0) and a follow-up measurement (y_1) after the 6-month treatment. Due to randomization, $E(y_0)$ is the same for all patients regardless of treatment assignment. Let $z = y_1 - y_0$. The null hypothesis is that $E(z)$ is the same for all participants regardless of treatment assignment. The data analysis is the same as for design option 1, but the outcome is z rather than y_1 . The approach is equivalent to fitting a general linear multivariate model with two outcome variables, y_0 and y_1 , and using $\mathbf{U} = [1 \ -1]'$.

Participants are statistically independent with data assumed to be Gaussian. Which design yields the better power depends on unknown parameters. The anticipated difference between follow-up and baseline, $E(y_1 - y_0)$, is 0.75 for Cimetidine, 1.00 for Ranitidine, and 1.25 for Famotidine.

Previous studies suggest that $\mathcal{V}(y_0) \approx \mathcal{V}(y_1) \approx 1$ and $\rho(y_0, y_1) \equiv \rho \approx 0.6$.

20.3.1 Under design 1, what is the probability of rejecting the null hypothesis when there are $n = 50$ participants per group ($N = 150$ participants total)?

20.3.2 What is the smallest sample size needed to achieve 0.80 power for design 1?

20.3.3 Under design 2, what is the probability of rejecting the null hypothesis when there are $n = 50$ participants per group ($N = 150$ participants total)?

20.3.4 What is the smallest sample size needed to achieve 0.80 power for design 2?

20.3.5 Repeat 20.3.3 using several different choices of ρ between 0 and 1. Plot a graph illustrating the relationship between ρ and power (an approximate, hand-drawn graph is acceptable).

20.3.6 The physician wants to enroll as many participants as possible, but can only spend a fixed amount of money. You're asked to help her determine which design is more cost effective. Under design option 1, the study costs \$75 per participant. Under design option 2, the study costs \$100 per participant. Which design do you recommend and why?

Sample Size for Multivariate Linear Models

21.1 THE MACHINERY OF A POWER ANALYSIS

We focus here mostly on the theory of power analysis and ignore the many important practical issues surrounding the task. The reader may wish to consult O'Brien and Muller (1993) for a tutorial and many related references.

The discussion assumes Gaussian errors and fixed predictors. For data analysis, the presence of random predictors causes no additional complexity. In contrast, allowing random predictors introduces an additional layer of distributional complexity to power analysis. Glueck and Muller (2003) reviewed the limited range of cases that have been solved. The same authors considered fixed predictors in combination with Gaussian covariates. An important open question remains concerning how to compute power for the interaction of fixed and Gaussian predictors. More generally, power analysis for random but not Gaussian predictors have received little attention.

We restrict attention to testable hypotheses, with full-rank \mathbf{C} , \mathbf{U} , Σ_* , and \mathbf{M} , while $\mathbf{C} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}(\mathbf{X}'\mathbf{X})$. As discussed in Section 2.5 of Muller, LaVange, Ramey, and Ramey (1992), α , Σ , \mathbf{X} , \mathbf{B} , \mathbf{C} , \mathbf{U} , and Θ_0 fully specifies a multivariate general linear model power calculation. Both α , the size of the test, as well as the test statistic must be chosen a priori. The dimensions of the model fix the degrees of freedom. Although specifying \mathbf{B} and Σ suffices to complete the power analysis, specifying the (usually smaller and simpler) matrices $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$, Θ_0 , and $\Sigma_* = \mathbf{U}'\Sigma\mathbf{U}$ also suffices and often proves easier. Further simplification occurs because the $b \times b$ matrix $\Omega = (\Theta - \Theta_0)' \mathbf{M}^{-1} (\Theta - \Theta_0) \Sigma_*^{-1}$ suffices. As a final simplification, the eigenvalues of Ω provide the minimally sufficient (additional) information required to complete the power analysis. Eigenvalue k of Ω is $\omega_k = N\rho_k^2 / (1 - \rho_k^2)$, with ρ_k^2 a (generalized) squared canonical correlation. Also $0 \leq \omega_k < \infty$ and $0 \leq \rho_k^2 \leq 1$. Writing $\Sigma_* = \Phi\Phi'$ gives $\Sigma_*^{-1} = \Phi^{-t}\Phi^{-1}$ and allows defining the symmetric matrix $\Omega_\Phi = \Phi^{-1}(\Theta - \Theta_0)' \mathbf{M}^{-1} (\Theta - \Theta_0) \Phi^{-t}$, which is a quadratic form in the symmetric matrix \mathbf{M}^{-1} . The eigenvalues of Ω equal the eigenvalues of Ω_Φ . For $(a \times b)$ Θ , $s = \min(a, b)$ and $\text{rank}(\Theta) = s_*$,

with $0 \leq s_* \leq s$. In turn, $\text{rank}(\Omega) = s_*$, the number of nonzero ω_k , which equals the number of nonzero ρ_k^2 . Also $s_* = 0 \Leftrightarrow \omega_k \equiv 0 \Leftrightarrow \rho_k^2 \equiv 0 \Leftrightarrow \Theta = \Theta_0 \Leftrightarrow H_0 \text{ holds} \Leftrightarrow \text{power} = \alpha$.

The sizes of s and s_* control many features of power analysis in the multivariate GLM. If $s = 1$, then all tests coincide and $s_* = 1$. Having $s > 1$ allows increasing s_* , which always increases power (with all other properties held constant). The theoretical property has a very important and practical implication. A study design with sufficient power to detect a linear dose effect will have even more power if an additional quadratic effect occurs (above and beyond the linear). Hence power analysis with $s_* = 1$ often may be taken as a conservative approach.

Deleting any duplicate rows from the design matrix (\mathbf{X}) creates the essence matrix (Definition 11.5). It simplifies determining properties of a design, such as rank. Comparing essence matrices allows determining relationships between alternate parameter sets. The concept provides convenient separation of total sample size from the coding scheme. The separation helps simplify interpretation and computation of power.

21.2 PAIRED t TEST EXAMPLE

A paired data t test can be conducted with $\text{Es}(\mathbf{X}) = \mathbf{I}_1 = 1$, $\mathbf{C} = 1$, and $\mathbf{U} = [-1 \ 1]'$. The power for $H_0 : \mu_2 = \mu_1$ will be the same for either of two different models:

$$\mathbf{Y} = \mathbf{1}_n[\mu_1 \ \mu_2] + \mathbf{E} \quad (21.1)$$

$$\mathbf{Y} = \mathbf{1}_n[0 \ 1]\delta + \mathbf{E}. \quad (21.2)$$

Here $\mathbf{C}\mathbf{B}\mathbf{U} = \theta = (\mu_2 - \mu_1) = \delta$. The second model avoids the need to specify the grand mean, $(\mu_2 + \mu_1)/2$, which has no effect on the power analysis.

Example 21.1 The code can use $\text{BETA}=\{0 \ 1\}=\{(\{0,1\})'$ and $\text{BETASCAL}=\delta$ or $\text{BETASCAL}=\text{DO}(\delta_{\text{low}}, \delta_{\text{high}}, \delta_{\text{increment}})$;

```
TITLE1 "P0803.SAS--simple paired t example";
PROC IML SYMSIZE=4000 WORKSIZE=4000;
%INCLUDE "..\IML\POWERLIB.IML";
U=({1 -1})`;
C={1}; THETA0=0;
ALPHA=.01;
SIGSCAL={1};
RHOSCAL={1};
SIGMA={ 2.1 3.2 ,
        3.2 2.4 };
*or; P=2; VARIANCE=2.1; RHO=.4;
      SIGMA=VARIANCE#( I(P)#(1-RHO) +J(P,P,RHO) );
ESSENCEX=I(1);
REPN={5,10};
```

```
BETA={0 1};
BETASCAL=DO(0, .30, .15);
RUN POWER;
```

__HOLDPOW	CASE	ALPHA	SIGSCAL	RHOSCAL	BETASCAL	TOTAL_N	POWER
	1	0.01	1	1	0	5	0.01
	2	0.01	1	1	0	10	0.01
	3	0.01	1	1	0.15	5	0.011
	4	0.01	1	1	0.15	10	0.012
	5	0.01	1	1	0.3	5	0.013
	6	0.01	1	1	0.3	10	0.02

21.3 TIME BY TREATMENT EXAMPLE

For $GLM_{N,p,q}(Y_i; X_i B, \Sigma)$ with Gaussian errors, the hypothesis of time-by-treatment interaction often generates particular interest. All participants are assumed to have been measured at the same times, $\{t_1, \dots, t_p\}$, with $Es(X) = I_q$ and $X = \mathbf{1}_n \otimes Es(X)$. In the simplest case $s_* = 1$.

The hypothesis of treatment-by-time interaction may be expressed in terms of differences of differences of means, with

$$C_1 = [\mathbf{1}_{q-1} \quad -I_{q-1}] \tag{21.3}$$

$$U_1 = [\mathbf{1}_{p-1} \quad -I_{p-1}]' \tag{21.4}$$

$$\Theta_0 = \mathbf{0}. \tag{21.5}$$

Useful corresponding canonical forms are, with $q \times p$ B_1 and $a \times b$ Θ_1 ,

$$\begin{aligned}
 B_1 &= \begin{bmatrix} \theta & \mathbf{0}_{1 \times (p-1)} \\ \mathbf{0}_{(q-1) \times 1} & \mathbf{0}_{(q-1) \times (p-1)} \end{bmatrix} \\
 &= \begin{bmatrix} 1 & \mathbf{0}_{1 \times (p-1)} \\ \mathbf{0}_{(q-1) \times 1} & \mathbf{0}_{(q-1) \times (p-1)} \end{bmatrix} \delta, \tag{21.6}
 \end{aligned}$$

$$\begin{aligned}
 \Theta_1 &= \begin{bmatrix} \theta & \mathbf{0}_{1 \times (b-1)} \\ \mathbf{0}_{(a-1) \times 1} & \mathbf{0}_{(a-1) \times (b-1)} \end{bmatrix} \\
 &= \begin{bmatrix} 1 & \mathbf{0}_{1 \times (b-1)} \\ \mathbf{0}_{(a-1) \times 1} & \mathbf{0}_{(a-1) \times (b-1)} \end{bmatrix} \delta. \tag{21.7}
 \end{aligned}$$

The alternative hypothesis of interest may be a linear-by-linear (treatment-by-time) interaction. If so, C_2 and U_2 give the orthonormal polynomial trends for q and p values, respectively. If $p = 3$ and $q = 4$ (assuming equal spacing for both between and within factors), $\Theta_0 = \mathbf{0}$ and

$$C_2 = \text{Dg}(\{20, 4, 4\})^{-1/2} \begin{bmatrix} -3 & -1 & 1 & 3 \\ 1 & -1 & -1 & 1 \\ 1 & -1 & 1 & -1 \end{bmatrix} = \begin{bmatrix} C_{2,1} \\ \vdots \\ C_{2,q-1} \end{bmatrix} \tag{21.8}$$

$$U_2 = \begin{bmatrix} -1 & 1 \\ 0 & -2 \\ 1 & 1 \end{bmatrix} \text{Dg}(2, 6)^{-1/2} = [u_{2,1} \ u_{2,2} \ \cdots \ u_{2,p-1}]. \tag{21.9}$$

Also, with $q \times p$ B_2 , $q \times 1$ $C'_{2,1}$, $1 \times p$ $u'_{2,1}$, $a \times b$ Θ_2 ,

$$B_2 = (C'_{2,1} u'_{2,1}) \delta \tag{21.10}$$

and

$$\begin{aligned} \Theta_2 &= C_2 (C'_{2,1} u'_{2,1} \delta) [u_1 \ u_1 \ \cdots \ u_{p-1}] \delta \\ &= \begin{bmatrix} C_{2,1} \\ \vdots \\ C_{2,q-1} \end{bmatrix} C'_{2,1} u'_{2,1} [u_1 \ u_1 \ \cdots \ u_{p-1}] \delta \\ &= \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} [1 \ 0 \ \cdots \ 0] \delta \\ &= \begin{bmatrix} 1 & \mathbf{0}_{1 \times (b-1)} \\ \mathbf{0}_{(a-1) \times 1} & \mathbf{0}_{(a-1) \times (b-1)} \end{bmatrix} \delta = \Theta_1. \end{aligned} \tag{21.11}$$

More generally, if $\delta = [\delta_1 \ \delta_2 \ \cdots \ \delta_{s_*}]'$, C_{2*} indicates the first s_* rows of C_2 , and U_{2*} indicates the first s_* columns of U_2 , then $B_{2*} = C'_{2*} \text{Dg}(\delta) U_{2*}$ has rank s_* and

$$\begin{aligned} \Theta_{2*} &= C_2 B_{2*} U_2 \\ &= (C_2 C'_{2*}) \text{Dg}(\delta) (U_{2*}' U_2) \\ &\quad [_{(a \times q)(q \times s_*)}][_{(s_* \times s_*)}][_{(s_* \times p)(p \times b)}] \\ &= \begin{bmatrix} I_{s_*} \\ \mathbf{0}_{(a-s_*) \times s_*} \end{bmatrix} \text{Dg}(\delta) [I_{s_*} \ \mathbf{0}_{s_* \times (b-s_*)}] \\ &= \begin{bmatrix} \text{Dg}(\delta) & \mathbf{0} \\ \mathbf{0} & \mathbf{0}_{(a-s_*) \times (b-s_*)} \end{bmatrix}. \end{aligned} \tag{21.12}$$

Choosing

$$B_3 = \Theta_{2*} \tag{21.13}$$

$$C_3 = \begin{bmatrix} I_{s_*} & \mathbf{0} \\ \mathbf{0}_{(a-s_*) \times s_*} & \mathbf{0} \end{bmatrix} \tag{21.14}$$

$$U_3 = \begin{bmatrix} I_{s_*} \\ \mathbf{0}_{(b-s_*) \times s_*} \end{bmatrix} \tag{21.15}$$

implies $C_3 B_3 U_3 = \Theta_{2*}$.

In some cases, the alternative hypothesis of interest may concern a main effect. For a between-groups main effect, examples are

$$B = \frac{1}{6} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 1 & 1 \\ 2 & 2 & 2 \\ 3 & 3 & 3 \end{bmatrix} \delta \tag{21.16}$$

and

$$B = C_{2,1} \mathbf{1}'_p \delta \tag{21.17}$$

while a within-subject main effect example is

$$B = \mathbf{1}'_q \mathbf{u}'_{2,1} \delta. \tag{21.18}$$

Example 21.2 The following code provides an example calculation of $C_2 C'_{2*}$.

```
PROC IML; RESET PRINT;
Q=4; SSTAR=2;
POLY=ORPOL(1:Q);
C=(POLY[* , 2:Q])`;
CSTAR=C[1:SSTAR, *];
CCSTARP=C*CSTAR`;
The following program computes power for a time by treatment
interaction.
PROC IML SYMSIZE=4000 WORKSIZE=4000;
%INCLUDE "..\IML\POWERLIB.IML";
P=3; U=( J(P-1, 1, 1) || (-I(P-1)) )`;
Q=4; C=J(Q-1, 1, 1) || (-I(Q-1));
ALPHA=.01;

VARIANCE=2.1;
RHO=.4;
SIGMA=VARIANCE#( I(P) # (1-RHO) + J(P, P, RHO) );
SIGSCAL={1, 2};
RHOSCAL={1};

ESSENCEX=I(Q);
REPN={5, 10};
BETA=J(Q, P, 0);
BETA[1, 1]=1;
BETASCAL=DO(0, .30, .15);

RUN POWER;
```

21.4 COMPARING BETWEEN AND WITHIN DESIGNS

A multivariate GLM with Gaussian errors for a paired data t test, with \mathbf{y}_j dimension $N \times 1$, may be written

$$[\mathbf{y}_1 \ \mathbf{y}_2] = \mathbf{1}_N [\mu_1 \ \mu_2] + \mathbf{E}. \quad (21.19)$$

For simplicity of comparison to an independent t test analysis, we assume homogeneity of variance holds across repeated measures (within subject):

$$\Sigma = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}. \quad (21.20)$$

Using $\mathbf{C} = [1]$ and $\mathbf{U} = [1 \ -1]'$ tests $H_0 : \mu_1 = \mu_2$. In turn, $\Theta = (\mu_1 - \mu_2)$,

$$\begin{aligned} \Sigma_* &= \mathbf{U}' \Sigma \mathbf{U} \\ &= \sigma^2 [1 \ -1] \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix} \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ &= \sigma^2 2(1 - \rho), \end{aligned} \quad (21.21)$$

$$\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{C}' = N^{-1}, \quad (21.22)$$

and

$$\begin{aligned} \Omega &= (\Theta - \Theta_0)' \mathbf{M}^{-1} (\Theta - \Theta_0) \Sigma_*^{-1} \\ &= \frac{(\mu_1 - \mu_2)^2}{\sigma^2} \left[\frac{N}{2(1 - \rho)} \right]. \end{aligned} \quad (21.23)$$

For $\nu_e = N - \text{rank}(\mathbf{X})$, the degrees of freedom are $ab = 1$ and $s(\nu_e - b + s) = N - 1$.

An independent groups t test with a balanced design corresponds to

$$\begin{bmatrix} \mathbf{y}_1 \\ \mathbf{y}_2 \end{bmatrix} = \mathbf{I}_2 \otimes \mathbf{1}_{N/2} \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix} + \mathbf{e}. \quad (21.24)$$

Assuming $\mathbf{e} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_N)$ implies homogeneity of variance (between groups). Choosing $\mathbf{C} = [1 \ -1]$ and $\mathbf{U} = [1]$ allows testing $H_0 : \mu_1 = \mu_2$. In turn, $\Theta = (\mu_1 - \mu_2)$, $\Sigma_* = \sigma^2$,

$$\mathbf{M} = [1 \ -1] [(N/2)\mathbf{I}_2]^{-1} \begin{bmatrix} 1 \\ -1 \end{bmatrix} = 4/N, \quad (21.25)$$

and

$$\Omega = \frac{(\mu_1 - \mu_2)^2}{\sigma^2} \left(\frac{N}{4} \right). \quad (21.26)$$

The degrees of freedom are $ab = 1$ and $s(\nu_e - b + s) = N - 2$.

Two possible comparisons seem interesting, based on having either a constant number of independent sampling units (ISUs) or a constant number of observations. With constant ISUs, N has the same value in the two settings, and the error degrees of freedom for paired minus independent equals $(N - 1) - (N - 2) = 1$. In turn,

$$\begin{aligned} \frac{\omega_{\text{paired}}}{\omega_{\text{indep}}} &= \frac{N/[2(1 - \rho)]}{N/4} \\ &= \frac{2}{(1 - \rho)} > 1. \end{aligned} \tag{21.27}$$

Both discrepancies favor the paired design (with constant ISUs, which requires collecting more observations).

With a constant number of observations, the error degrees of freedom for paired minus independent equals $(N/2 - 1) - (N - 2) = -N/2 + 1$. In turn,

$$\begin{aligned} \frac{\omega_{\text{paired}}}{\omega_{\text{indep}}} &= \frac{(N/2)/[2(1 - \rho)]}{N/4} \\ &= \frac{1}{(1 - \rho)}. \end{aligned} \tag{21.28}$$

If $\rho \leq 0$, then $(1 - \rho)^{-1} \leq 1$ and both discrepancies favor the independent design. If $\rho \geq 0$, then $(1 - \rho)^{-1} \geq 1$ and for each N there exists $\rho_0(N)$ such that $\rho \geq \rho_0$ implies the paired design is superior, and otherwise the independent design is superior. Equivalently, for each ρ there exists $N_0(\rho)$ such that $N \geq N_0(\rho)$ and the paired design is superior, while otherwise the independent design is superior. Some regions may be undefined. A three-dimensional plot, with power difference vertically, N and ρ as the floor plane, can be especially informative. Critical values are $f_p = F_F^{-1}(1 - \alpha; 1, N/2 - 1)$ and $f_1 = F_F^{-1}(1 - \alpha, 1, N - 2)$. In turn, for method m , $\text{Power}(m) = 1 - F_F(f_m; 1, \nu_m; \omega_m)$.

Other comparisons are also interesting. For paired data

$$\begin{aligned} \text{vec}(\hat{\Theta}) &\sim \mathcal{N}_{ab}[\text{vec}(\Theta), \Sigma_* \otimes \mathbf{M}] \\ &\sim \mathcal{N}[(\mu_1 - \mu_2), 2(1 - \rho)\sigma^2/N], \end{aligned} \tag{21.29}$$

while for independent data

$$\begin{aligned} \text{vec}(\hat{\Theta}) &\sim \mathcal{N}_{ab}[\text{vec}(\Theta), \Sigma_* \otimes \mathbf{M}] \\ &\sim \mathcal{N}[(\mu_1 - \mu_2), 4\sigma^2/N]. \end{aligned} \tag{21.30}$$

With equal sample sizes

$$\begin{aligned} \gamma &= \frac{2(1 - \rho)\sigma^2/N}{4\sigma^2/N} \\ &= (1 - \rho)/2, \end{aligned} \tag{21.31}$$

while for equal observations

$$\begin{aligned}\gamma &= \frac{2(1 - \rho)\sigma^2/(N/2)}{4\sigma^2/N} \\ &= (1 - \rho)/4.\end{aligned}\tag{21.32}$$

The difference in error degrees of freedom affects the confidence interval. Hence one should compare quantiles of confidence interval widths, which are of the form

$$w = 2\sigma_m \sqrt{F_F^{-1}(1 - \alpha; 1, \nu_m)}.\tag{21.33}$$

21.5 SOME INVARIANCE PROPERTIES

Any testable secondary parameter $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$ ($a \times b$), with $H_0 : \Theta = \Theta_0$, has $\text{rank}(\mathbf{C}) = a \leq q$ and $\text{rank}(\mathbf{U}) = b \leq p$. Eigenvalues of $\mathbf{S}_h\mathbf{S}_e^{-1}$, and hence all multivariate test statistics and associate p values, are invariant to full-rank transformation of rows of \mathbf{C} and columns of \mathbf{U} (Section 16.8). Specifying \mathbf{B} and Σ usually requires the most thought. Specifying $\Theta = \mathbf{C}\mathbf{B}\mathbf{U}$, Θ_0 , and $\Sigma_* = \mathbf{U}'\Sigma\mathbf{U}$ suffices. In turn, \mathbf{B} and Σ reduce to canonical forms. For $H_0 : \mu_1 = \mu_2$ either $\beta = [\mu_1 \mu_2]'$ or $\beta = [0 \delta]'$ (with $\delta = \mu_2 - \mu_1$) leads to the same power. In practice, most tests involve hypotheses which exclude the intercept, while the model does span an intercept. Such situations allow assuming a grand mean of zero because the mean has no effect on power.

For fixed α , Σ , \mathbf{X} , \mathbf{C} , \mathbf{U} , and Θ_0 , power for any $s_* = 1$ alternative may be expressed and plotted as a function of a scalar parameter, such as a mean difference or a squared canonical correlation. Choosing such a representation can be extremely helpful in defining a range of alternate hypotheses of interest. Plotting power in terms of the scalar parameter is especially enlightening.

21.6 RANDOM PREDICTORS

As discussed in the previous chapter, the presence of random predictors greatly complicates noncentral distribution theory. Sampson (1974) detailed many of the issues for the univariate and multivariate models with Gaussian predictors.

Glueck and Muller (1998) reviewed the limited work on random predictor power in multivariate models. They also described methods for accurately approximating power for a $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i; \mathbf{B}, \Sigma)$ with combinations of fixed and Gaussian predictors corresponding to a baseline covariate design. Approximating power for the more general model allowing unequal slopes for each group remains an open and important question.

21.7 INTERNAL PILOT DESIGNS

The principles and advantages of internal pilot designs for univariate models were introduced in Section 20.7. An interim residual analysis (without any interim data analysis) provides a variance estimate $\hat{\sigma}_1^2$. In turn, a power analysis based on $\hat{\sigma}_1^2$ leads to increasing or decreasing total sample size. Avoiding the potential for inflated test size due to using an internal pilot requires special testing procedures.

Coffey and Muller (2003) considered using internal pilot designs with the UNIREP approach to repeated measures. They (incorrectly) speculated that the inherent conservatism of the Geisser-Greenhouse test in small samples might compensate for the test size inflation induced by an internal pilot. The interim analysis produces $\hat{\Sigma}_{*1} = \mathbf{U}'\hat{\Sigma}_1\mathbf{U}$ and an updated sample size choice. Simulations demonstrated that test size was inflated above the target level in small samples. Hence work on developing other strategies for controlling test size was begun.

Coffey and Muller drew an additional conclusion from their simulations. For $\epsilon = \text{tr}^2(\Sigma_*) / [\text{btr}(\Sigma_*^2)]$ near the lower boundary of $1/b$ and small N , the UNIREP power approximations of Muller and Barton (1989) lack sufficient accuracy for internal pilot use. In the worst case, $b = 1$, $\epsilon \approx 0.29$, and $N = 20$, a particular pattern of means gave a predicted power of 0.65, while a power of 0.87 was observed in a simulation with 100,000 replications (standard error ≈ 0.0015). Muller, Edwards, Simpson, and Taylor (2006) demonstrated approximations for power of the UNIREP tests which almost entirely eliminate the inaccuracy. In the worst-case condition just described, the new method predicted power of 0.84.

EXERCISES

21.1 Use the POWERLIB software described in the Appendix A (Section A.2) to reproduce the results in Example 21.1.

21.2 A clinical trial is planned to compare a new drug with an existing one. Treatment starts on a Monday morning. The outcome is measured that afternoon and the four following afternoons. A multivariate $\text{GLM}_{N,p,q}(\mathbf{Y}_i; \mathbf{X}_i\mathbf{B}, \Sigma)$ with Gaussian errors has \mathbf{Y} (20×5), and

$$\mathbf{XB} = \begin{bmatrix} \mathbf{1}_{10} & \mathbf{0}_{10} \\ \mathbf{0}_{10} & \mathbf{1}_{10} \end{bmatrix} \begin{bmatrix} \mu_{11} & \mu_{12} & \mu_{13} & \mu_{14} & \mu_{15} \\ \mu_{21} & \mu_{22} & \mu_{23} & \mu_{24} & \mu_{25} \end{bmatrix}.$$

Suppose that $\Sigma = \sigma^2[(1 - \rho)\mathbf{I} + \mathbf{1}\mathbf{1}'\rho]$, $\sigma^2 = 7$, $\rho = 0.5$, and $\mu_{ij} = (i - 1)(j - 1)$ for $i \in \{1, 2\}$, $j \in \{1, \dots, 5\}$. The Wilks statistic has been chosen for all tests. Define

$$\Theta_1 = [\mu_{11} - \mu_{21} \quad \mu_{12} - \mu_{22} \quad \mu_{13} - \mu_{23} \quad \mu_{14} - \mu_{24} \quad \mu_{15} - \mu_{25}]$$

$$\theta_2 = [(\mu_{11} + \mu_{12} + \mu_{13} + \mu_{14} + \mu_{15})/5 - (\mu_{21} + \mu_{22} + \mu_{23} + \mu_{24} + \mu_{25})/5]$$

$$\Theta_3 = \begin{bmatrix} \mu_{11} - \mu_{12} & \mu_{11} - \mu_{13} & \mu_{11} - \mu_{14} & \mu_{11} - \mu_{15} \\ \mu_{21} - \mu_{22} & \mu_{21} - \mu_{23} & \mu_{21} - \mu_{24} & \mu_{21} - \mu_{25} \end{bmatrix}$$

$$\Theta_4 = (\mu_{11} + \mu_{21}) \cdot \mathbf{1}'_4 - [(\mu_{12} + \mu_{22}) \quad (\mu_{13} + \mu_{23}) \quad (\mu_{14} + \mu_{24}) \quad (\mu_{15} + \mu_{25})].$$

21.2.1 Write a one-sentence scientific interpretation, aimed at the scientists, of $H_0 : \Theta_1 = \mathbf{0}$.

21.2.2 Specify the \mathbf{C} and \mathbf{U} matrices needed.

21.2.3 What is the probability of rejecting $H_0 : \Theta_1 = \mathbf{0}$ using a 0.05 level test?

21.2.4 Write a one-sentence scientific interpretation, aimed at the scientists, of $H_0 : \theta_2 = 0$.

21.2.5 Specify the \mathbf{C} and \mathbf{U} matrices needed.

21.2.6 What is the probability of rejecting the null hypothesis $\theta_2 = 0$ using a 0.05 level test?

21.2.7 Write a one sentence scientific interpretation, aimed at the scientists, of $H_0 : \Theta_3 = \mathbf{0}$.

21.2.8 Specify the \mathbf{C} and \mathbf{U} matrices needed.

21.2.9 What is the probability of rejecting the null hypothesis $\Theta_3 = \mathbf{0}$ using a 0.05 level test?

21.2.10 Write a one-sentence scientific interpretation, aimed at the scientists, of $H_0 : \Theta_4 = \mathbf{0}$.

21.2.11 Specify the \mathbf{C} and \mathbf{U} matrices needed.

21.2.12 What is the probability of rejecting the null hypothesis $\Theta_4 = \mathbf{0}$ using a 0.05 level test?

21.2.13 The various tests arise from either a MANOVA or repeated-measures approach. A third alternative is a set of Bonferroni corrected tests of drug difference, one for each day.

21.2.14 Specify the \mathbf{C} and \mathbf{U} matrices needed to conduct the five tests and the appropriate nominal level of α .

21.2.15 What is the probability of rejecting the null hypothesis $H_0 : \Theta_{5,1} = 0$?

21.2.16 What is the probability of rejecting the null hypothesis $H_0 : \Theta_{5,5} = 0$?

21.2.17 In practice, balancing control of test size and maximizing power leads to using only one of the approaches (MANOVA, REPM, Bonferroni univariate). Which seems preferable? Which particular tests are most logically consistent (not necessarily the most powerful) with the speculation that $\mu_{ij} = (i-1)(j-1)$? Without doing any additional calculations, do the results suggest a different sample size, assuming a target power of 0.90 or better?

21.3 Consider Table 6, p. 553, in Muller and Barton (1989). Row 9 provides data about condition 113, which is detailed in their Table 4 and associated text.

21.3.1 The canonical form of the \mathbf{B} matrix is described in the right-hand column of p. 552, line 9. Think of the form as $\mathbf{B} = c\mathbf{B}_c$. For the conditions of the simulations associated with their Tables 4–6, give a numerical value for \mathbf{B}_c . Write

this as a matrix of numbers.

21.3.2 For the conditions of the simulations associated with their Tables 4–6, use their Table 4 and associated text to give a numerical value of the covariance matrix used for condition 113. Write the results as a matrix of numbers.

21.3.3 Use power software to find the value of c which was used to produce the predicted power value of 0.80 for the GG test for the ninth line of Table 6 (p. 553), for condition 113.

Hint: In SAS/IML[®], use the DO function to create a list of candidate c values and assign the result to BETASCAL, as in the examples.

21.3.4 Produce a high-resolution plot of power as a function of c .

Hint: You will likely wish to use the DS option in POWERLIB.

21.4 Assume the notation and setting of Muller, LaVange, Ramey, and Ramey, (1992). As throughout, $\mathbf{1}_j$ indicates a $j \times 1$ vector of 1's. With $n = 9$, define $\mathbf{X} = \mathbf{I}_3 \otimes \mathbf{1}_n$,

$$\mathbf{B} = \mu \mathbf{1}_3 \mathbf{1}'_5 + \delta \frac{1}{16} \begin{bmatrix} 0 & 1 & 4 & 9 & 16 \\ 0 & 1 & 2 & 3 & 4 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix},$$

$$\mathbf{\Sigma} = \sigma^2 \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 & 5 \end{bmatrix}^{1/2} \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 & \rho^4 \\ \rho & 1 & \rho & \rho^2 & \rho^3 \\ \rho^2 & \rho & 1 & \rho & \rho^2 \\ \rho^3 & \rho^2 & \rho & 1 & \rho \\ \rho^4 & \rho^3 & \rho^2 & \rho & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 3 & 0 & 0 \\ 0 & 0 & 0 & 4 & 0 \\ 0 & 0 & 0 & 0 & 5 \end{bmatrix}^{1/2},$$

$$\mathbf{C} = \begin{bmatrix} 1 & -1 & 0 \\ 1 & 0 & -1 \end{bmatrix},$$

$$\mathbf{U}' = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ 1 & 0 & -1 & 0 & 0 \\ 1 & 0 & 0 & -1 & 0 \\ 1 & 0 & 0 & 0 & -1 \end{bmatrix},$$

and $\mathbf{\Theta}_0 = \mathbf{0}$.

Assume $\delta = 2$, $\sigma^2 = 1$, and $\rho = 0.6$. Use $\alpha = 0.025$.

21.4.1 Prove that $\mathbf{\Theta}$ is invariant to the value of μ for this hypothesis. You may use IML, or any other matrix language, for the purely numerical calculations. (Hence, without loss of generality, assume $\mu = 0$.)

21.4.2 Compute $\mathbf{\Omega}$.

21.4.3 Compute the eigenvalues of $\mathbf{\Omega}$.

21.4.4 Compute the squared canonical correlations.

21.4.5 Compute the (population values) of the measures of multivariate association for the four MULTIREP (invariant) tests.

21.5 Write your own matrix language code to directly implement the steps in Section 2.5 of Muller, LaVange, Ramey, and Ramey (1992, p. 1214) to compute

the approximate power of the Pillai-Bartlett test *only*. In particular, compute

21.5.1 Step 1 in Section 2.5 for the Pillai-Bartlett test only.

21.5.2 Step 2 in Section 2.5 for the Pillai-Bartlett test only.

21.5.3 Step 3 in Section 2.5 for the Pillai-Bartlett test only.

21.5.4 Step 4 in Section 2.5 for the Pillai-Bartlett test only.

21.6. Use POWERLIB to compute the same power value as in the last part of the exercise.

Sample Size for Generalizations of Multivariate Models

22.1 MOTIVATION

In developing new statistical methods, statisticians have historically focused first on estimation and then on inference. Methods for choosing a sample size typically come last. Practical and mathematical reasons stimulate the order. Distribution theory for sample size involves greater complexity than estimation or inference under the null. Furthermore, the enthusiasm that practicing data analysts have for sample size methods has never been shared by more theoretical statisticians. Hence few results are available for power and sample size analysis of generalizations of multivariate linear models. In keeping with our discussions of estimation and testing, we sketch some results for growth curve models.

22.2 SAMPLE SIZE METHODS FOR GROWTH CURVE MODELS

Some noncentral theory has been developed for growth curves with higher order trends used as covariates, especially for large samples. However, the methods have not been studied carefully in small samples. Berger (1986) used simulations to study both test size and power for a variety of methods for analyzing growth curves. His results support avoiding the use of high-order polynomials as covariates due to the likelihood of inflating test size.

Without the use of high-order polynomials as covariates, growth curve analysis reduces to a special case of a multivariate GLM. Hence methods in the previous chapter apply. Particular care must be taken to choose U matrices to correctly reflect the specified model.

An example will illustrate the issue. If the data involve five time points while the desired model includes only linear and quadratic trends (along with zero order), then two within-subject contrasts would be used, $U'_0 = 5^{-1/2} \cdot [1 \ 1 \ 1 \ 1 \ 1]$ and

$$\mathbf{U}'_{1-2} = \begin{bmatrix} 10 & 0 \\ 0 & 14 \end{bmatrix}^{-1/2} \begin{bmatrix} -2 & -1 & 0 & 1 & 2 \\ 2 & -1 & -2 & -1 & 2 \end{bmatrix}. \quad (22.1)$$

The omitted cubic and quartic trends are spanned by

$$\mathbf{U}'_{3-4} = \begin{bmatrix} 10 & 0 \\ 0 & 70 \end{bmatrix}^{-1/2} \begin{bmatrix} -1 & -2 & 0 & -2 & 1 \\ 1 & -4 & 6 & -4 & 1 \end{bmatrix}. \quad (22.2)$$

The default choices for a multivariate linear model (the multivariate approach to repeated measures, MULTIREP) would use \mathbf{U}_0 and

$$\mathbf{U}'_W = \begin{bmatrix} -1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 1 & 0 & 0 \\ -1 & 0 & 0 & 1 & 0 \\ -1 & 0 & 0 & 0 & 1 \end{bmatrix}. \quad (22.3)$$

However, the columns of \mathbf{U}_W span \mathbf{U}_{1-2} and \mathbf{U}_{3-4} . More precisely, if $\mathbf{U}_T = [\mathbf{U}_{1-2} \ \mathbf{U}_{3-4}]$ then $\mathbf{U}_T = \mathbf{U}_W \mathbf{P}$ for 4×4 and full rank. Invariance properties of the multivariate model imply that any test with \mathbf{U}_W gives the same result as a test with $\mathbf{U}_W \neq \mathbf{U}_{1-2}$.

Using \mathbf{U}_W rather than \mathbf{U}_{1-2} provides another example of an alignment error, as discussed in the previous chapter. Checking the degrees of freedom associated with the hypothesis and parameters leads to recognizing that $4 \neq 2$. Making sure dimensions and degrees of freedom correspond to the desired inference can help avoid alignment errors.

Sample Size for Linear Mixed Models

23.1 MOTIVATION

Mixed models have become one of the most widely used methods for data analysis. Unfortunately, as Verbeke and Molenberghs (2000, Section 23.2) noted, very little is known about nonnull distributions in mixed models. The wide variety of test statistics used in mixed models adds a substantial complexity to the task. Additional complexity arises from interest in random predictors. As always with questions of power, Monte Carlo simulations provide a completely defensible, although onerous, method of approximating power.

In linear models, full specification of noncentral distributions requires knowing the (population) distributions of any random predictors. For mixed models, the comment applies to both \mathbf{X}_i and \mathbf{Z}_i .

23.2 METHODS

Two approaches have been suggested for approximating power in general mixed models. Both are based on a supposition that reflects a property of many univariate and multivariate linear models. Only limited simulations are available to support the methods which use F approximations for Wald-type tests. However, the basic ideas are promising. In the particular case of a $\text{GLM}_{N,q}(y_i; \mathbf{X}_i\boldsymbol{\beta}, \sigma^2)$ with fixed predictors and Gaussian errors, the (testable) general linear hypothesis is $H_0: \mathbf{C}\boldsymbol{\beta} = \boldsymbol{\theta}_0$. With $\omega = (\boldsymbol{\theta} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1}(\boldsymbol{\theta} - \boldsymbol{\theta}_0)/\sigma^2$ and $\mathbf{M} = \mathbf{C}(\mathbf{X}'\mathbf{X})\mathbf{C}'$, the usual F statistic (for a *univariate* model) is given in equation 2.32 as $F = [(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)' \mathbf{M}^{-1}(\widehat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)/a]/\widehat{\sigma}^2 = \widehat{\omega}/a \sim F(a, N - r, \omega)$. Hence the noncentrality parameter may be characterized as $\omega = f_A a$, with f_A the F statistic occurring in the very special case with $\widehat{\boldsymbol{\theta}} = \boldsymbol{\theta}$ and $\widehat{\sigma}^2 = \sigma^2$. O'Brien and Muller (1993) described the concept as the *exemplary data* approach to power calculation and credited Graybill (1976) with earlier promotion of the idea. The exemplary data approach leads to noncentral F approximations for power of the mixed model.

Helms (1992) described a noncentral F power approximation for an F statistic (F_H) defined in terms of a novel modification of the usual REML estimators (equations 14–15 in his paper). Verbeke and Molenberghs (2000, Section 23.2)

provided some related discussion. However, the latter authors implicitly allow using standard REML or ML variance estimators (in the definition of the F statistic, their equation 6.6) while emphasizing the importance of specifying the degrees of freedom parameter. As mentioned earlier, additional simulations seem to be necessary to assess the accuracy of the approximations discussed.

Stroup (2003) described a strategy for approximating mixed model power which also uses the exemplary data concept. A simulation with 500 replications for one design gave power estimates consistent with the approximation. For computing convenience, Stroup used the Satterthwaite approach to determine approximate error degrees of freedom in power approximation. His simulations (for a narrow range of models) used the Kenward and Roger (1997) null case approximation (which currently provides the best control of test size).

23.3 INTERNAL PILOT DESIGNS

The principles and advantages of internal pilot designs for univariate models were introduced in Section 20.7. Zucker and Denne (2002) examined internal pilot designs for a two group clinical trial with at least $N = 40$ participants. In most cases, they disallowed a decrease in sample size. They used a likelihood ratio test with a second-order (Bartlett) approximation specifically derived for the covariance structure chosen, as derived by Zucker, Lieberman, and Manor (2000). Applying the method to any other design requires design-specific derivatives. An adjustment for using an internal pilot was required to avoid test size inflation.

APPENDIX

Computing Resources

A.1 EXAMPLE DATA: DETECTING BREAST CANCER IN MAGNETIC RESONANCE IMAGING WITH A CONTRAST AGENT

Research data used by permission of Dr. M. P. Braeuning and Dr. E. D. Pisano.

Approximately 1 in 9 American women develop breast cancer, based on full life expectancy. Breast cancer kills more American women than any other cancer. Regular self-examination for all women and screening mammography (x-rays of the breasts) for older women provide the main lines of defense. Having detected a suspicious region, the physician must discriminate between malignant and benign tissue. The life-and-death consequences of diagnosis make expensive procedures worthwhile. More information can be found at www.cancer.org.

The goal was to compare image brightness over time in different types of breast tissue on MRIs that followed injection of the contrast agent intravenous gadolinium-DTPA. For each woman, average brightness was recorded on each of four images, taken at 0, 45, 90, or 135 sec (+/- a few seconds) after injection. For each image, measurements were recorded from one region of fat, one of parenchyma, and one or two regions of interest (ROI). Each ROI was classified as benign or malignant (cancer), based on a subsequent pathologist's reading of a biopsy. Ten women requiring diagnosis were imaged. The following diagram shows the data available for each patient. Multicenter trials typically pay for a second pathology reading to increase reliability and validity of the (high quality but not perfect) single reading.

Patient	Cancer	Benign	Fat	Parenchyma
2	✓	✓	✓	✓
6	✓	✓	✓	✓
7	✓	✓	✓	✓
1	✓		✓	✓
3		✓	✓	✓
4	✓		✓	✓
5	✓		✓	✓
8		✓	✓	✓
9		✓	✓	✓
10	✓		✓	✓

Each ✓ indicates the presence of four observations (baseline=0, 45, 90, and 135 sec). Except where indicated in the diagram, for each patient, tissue type and time, the mean recorded MRI signal, the standard deviation, and the area (mm²) were available. Missing data indicates that no such tissue was found in the images.

Conversations with physicists led us to conclude that the logarithm of (signal at time j/signal at time 0) is proportional to concentration, which was chosen as the response variable. Exploratory analysis of Box-Cox power transformations of the response, coupled with evaluation of jackknife residuals, supported the choice.

All files can be downloaded from the Web site <http://ehpr.ufl.edu/muller/>. Raw data are in P0101.DAT (a text file). Programs P0101.SAS–P0105.SAS created SAS files P0101.SD2–P0105.SD2 (version 6.12, created on a PC). If you have trouble transporting the files to another platform, and feel compelled to start from raw data, you must use the programs provided to ensure that the same variable names, labels, and data are being analyzed.

The exercises use P0104.SD2, produced by P0104.SAS, which has the following statement:

```

LABEL DLOGROI1="Dif log ROI 45sec"
      DLOGROI2="Dif log ROI 90sec"
      DLOGROI3="Dif log ROI 135sec"
      DLOG_F_1="Dif log Fat 45sec"
      DLOG_F_2="Dif log Fat 90sec"
      DLOG_F_3="Dif log Fat 135sec"
      DLOG_P_1="Dif log Parenchyma 45sec"
      DLOG_P_2="Dif log Parenchyma 90sec"
      DLOG_P_3="Dif log Parenchyma 135sec"
      BENIGN  ="1=> Benign, else 0"
      MALIGN  ="1=> Malign, else 0";
    
```

The program P0104.SAS produced the following:

P0104.SAS>Create file using only one benign or malig per case
 All variables in file

	D	D	D	D	D	D	D	D	D	D	C
	L	L	L	L	L	L	L	L	L	L	O
	O	O	O	O	O	O	O	O	O	O	O
	G	G	G	G	G	G	G	G	G	G	G
	—	—	—	—	—	—	—	R	R	R	R
O	F	F	F	P	P	P	O	O	O	O	O
b	I	—	—	—	—	—	I	I	I	I	I
s	D	1	2	3	1	2	3	1	2	3	T
											N
1	2	0.002	0.008	0.041	0.210	0.259	0.278	0.944	1.127	1.195	1 1 0
2	6	0.083	0.071	0.137	0.054	0.001	0.142	0.084	0.037	0.088	1 1 0
3	3	0.118	0.271	0.234	0.001	0.055	0.056	0.081	0.184	0.260	1 1 0
4	8	-.178	-.243	0.056	0.007	0.056	0.066	0.087	0.209	0.193	1 1 0
5	9	0.137	0.080	0.191	0.174	0.162	0.127	0.062	0.197	0.184	1 1 0
6	7	0.016	0.165	0.135	0.038	0.149	0.168	0.121	0.162	0.261	1 0 1

```

7  1  0.054  0.057  0.017  0.037  -.031  0.019  0.099  0.003  0.082  1  0  1
8  4  0.623  0.518  0.604  -.096  0.180  0.154  -0.578  -0.124  0.092  1  0  1
9  5  -.085  0.231  0.153  -.033  0.026  0.071  0.071  0.136  0.293  1  0  1
10 10  -.122  -.285  -.084  0.029  0.070  0.064  0.232  0.225  0.253  1  0  1

```

Before you begin the exercises, please read programs P0101.SAS through P0104.SAS, as well as the associated log and list files (contained in the same directory) to understand both the programming and scientific decisions behind creating the file P0104.SD2. Most importantly, note that P0104.SD2 and P0105.SD2 have no missing data, while P0101.SD2–P0103.SD2 do.

A.2 FREE SOFTWARE

A.2.1 Overview

The software described here is available at no cost on the Web at <http://ehpr.ufl.edu>. All of the software is written in SAS/IML[®]. Hence the source code is included, which allows embedding the software in other programs, as well as translating modules to other languages. Each must be downloaded separately. The files include user manuals and many examples.

A.2.2 LINMOD: Multivariate Linear Models Analysis

LINMOD (LINear MODels) performs a wide variety of computations in SAS/IML for a general linear multivariate model with Gaussian errors. LINMOD allows the analyst familiar with matrix algebra notation and IML syntax to efficiently compute tests, estimates, and all associated statistics. Muller, LaVange, Ramey, and Ramey (1992) presented a succinct statement of the model, the general linear hypothesis, and associated statistics. Chapters 3 and 6 provide a more detailed overview. Chapters 4, 12, 13, 16, 17, 19, and 21 contain more detail.

All of the results can be computed with some combination of PROC GLM and REG in SAS. The primary advantage of LINMOD lies in the efficiency which results from the direct relationship between the syntax of the program and the matrix algebra formulation of models and tests. The primary disadvantage of LINMOD lies in the statistical sophistication required to use the program to produce valid results. However, anyone who has spent time trying to deduce exactly how an estimate, test, confidence interval, or correlation was computed by a particular option in a particular program will welcome the ability to explicitly control the calculations in a formula-based syntax. The additional overhead of LINMOD will usually only be worthwhile for complex designs or tests. The less simple and traditional the design and/or hypothesis, the more likely that LINMOD will appeal to someone able to define and analyze linear models in matrix notation. A further advantage lies in the fact that all of the computed values remain available to the user in matrices. This allows storing all results in a permanent SAS database (as distinct from the listing file) and exercising complete control of formatting.

Having results in matrix format makes subsequent custom processing very easy. Using ODS with GLM or REG provides an alternative approach.

Comparing LINMOD to PROC REG (or GLM) demonstrates a classic tradeoff in program design: increased flexibility and control for the sophisticated user versus friendly interface for the naive user. For example, LINMOD does not add or in any way recognize an intercept term in any model. The user must code one if desired and may choose to test it if present. Using LINMOD requires the ability to define and analyze linear models in matrix notation.

Starting from the assumption of the user having matrix knowledge, great effort was expended to give LINMOD an interface with consistent design, extensive error checking, and informative messages. The wide assortment of matrix operators and functions in PROC IML greatly enhances power and flexibility. Furthermore, the extensive editing, printing, plotting, and data management facilities of SAS are available for pre- and post-processing.

A.2.3 MISSMOD: Multivariate Linear Models with Missing Data

MISSMOD provides accurate test size in small samples for many kinds of Gaussian repeated measures and multivariate data with missing values. The software computes approximate tests described by Catellier and Muller (2000) for general linear multivariate models. Simulations support the conclusion that, in contrast to current mixed model competitors, the methods control test size even with as few as 12 observations for 6 repeated measures and 5% missing data. Assuming data missing at random (MAR), the EM algorithm provides maximum likelihood estimates using all of the available data. The tests generalize standard “multivariate” and “univariate” approaches to repeated-measures tests by reducing the error degrees of freedom by replacing the number of independent sampling units by various functions of the numbers of nonmissing pairs of responses. The program is based closely on LINMOD. Source code, an extensive user's guide, and example programs are included in the free download.

A.2.4 POWERLIB: Multivariate and Repeated-Measures Power

POWERLIB provides convenient power calculations for a wide range of multivariate linear models with Gaussian errors. The multivariate and univariate approaches to repeated measures as well as MANOVA tests are covered. F approximations are used throughout and are reduced to exact forms whenever possible. Approximate or exact power for the Wilks, Pillai-Bartlett, and Hotelling-Lawley tests is available. Approximate or exact power for the Box conservative test, Geisser-Greenhouse, Huynh-Feldt, and uncorrected test is also available. Confidence limits may be requested for most power values to reflect the uncertainty due to using estimated variances and, when appropriate, means and variances. A simple option causes SAS data files to be produced automatically, which simplifies producing plots and tables for manuscripts. Documentation

includes a range of designs and examples for both UNIX and Windows systems. Source code is included in the free download.

A.2.5 CISIZE: Sample Size Involving Confidence Intervals

CISIZE generalizes ideas about sample size to achieve confidence intervals and power properties. It computes $\Pr\{(W \cap R)|V\}$, $\Pr\{R\}$, $\Pr\{W\}$ and $\Pr\{W|V\}$ for any scalar hypothesis in a univariate or multivariate GLM with fixed predictors. The software implements the techniques discussed in Jiroutek et al. (2003).

A.3 ORTHOGONAL POLYNOMIAL COEFFICIENTS

$d =$ Number of Distinct Points, $p =$ Power of Trend

d	p	c_1	c_2	c_3	c_4	c_5	c_6	c_7	c_8	c_9	$\sum_{j=1}^d c_j^2$
2	1	-1	1	2
3	1	-1	0	1	2
	2	1	-2	1	6
4	1	-3	-1	1	3	20
	2	1	-1	-1	1	4
	3	-1	3	-3	1	20
5	1	-2	-1	0	1	2	10
	2	2	-1	-2	-1	2	14
	3	-1	2	0	-2	1	10
	4	1	-4	6	-4	1	70
6	1	-5	-3	-1	1	3	5	.	.	.	70
	2	5	-1	-4	-4	-1	5	.	.	.	84
	3	-5	7	4	-4	-7	5	.	.	.	180
	4	1	-3	2	2	-3	1	.	.	.	28
	5	-1	5	-10	10	-5	1	.	.	.	252
7	1	-3	-2	-1	0	1	2	3	.	.	28
	2	5	0	-3	-4	-3	0	5	.	.	84
	3	-1	1	1	0	-1	-1	1	.	.	6
	4	3	-7	1	6	1	-7	3	.	.	154
	5	-1	4	-5	0	5	-4	1	.	.	84
	6	1	-6	15	-20	15	-6	1	.	.	924
8	1	-7	-5	-3	-1	1	3	5	7	.	168
	2	7	1	-3	-5	-5	-3	1	7	.	168
	3	-7	5	7	3	-3	-7	-5	7	.	264
	4	7	-13	-3	9	9	-3	-13	7	.	616
	5	-7	23	-17	-15	15	17	-23	7	.	2184
	6	1	-5	9	-5	-5	9	-5	1	.	264
	7	-1	7	-21	35	-35	21	-7	1	.	3432
9	1	-4	-3	-2	-1	0	1	2	3	4	60
	2	28	7	-8	-17	-20	-17	-8	7	28	2772
	3	-14	7	13	9	0	-9	-13	-7	14	990
	4	14	-21	-11	9	18	9	-11	-21	14	2002
	5	-4	11	-4	-9	0	9	4	-11	4	468
	6	4	-17	22	1	-20	1	22	-17	4	1980
	7	-1	6	-14	14	0	-14	14	-6	1	858
	8	1	-8	28	-56	70	-56	28	-8	1	12870

References

- Anderson, T. W. (2004) *An Introduction to Multivariate Statistical Analysis*, 3rd ed., New York: Wiley.
- Andrade, D. F. and Helms, R. W. (1986) ML estimation and LR tests for the multivariate normal distribution with general linear model mean and linear-structure covariance matrix, *Communications in Statistics A*, **15**, 89–107.
- Andrews, D. F., Gnanadesikan, R. and Warner, J. L. (1971) Transformations of multivariate data, *Biometrics*, **27**, 825–840.
- Arnold, S. F. (1981) *The Theory of Linear Models and Multivariate Analysis*, New York: Wiley.
- Banerjee, K. S. (1964) A note on idempotent matrices, *Annals of Mathematical Statistics*, **35**, 880–882.
- Barton, C. N. and Cramer, E. C. (1989) Hypothesis testing in multivariate linear models with randomly missing data, *Communications in Statistics B*, **18**, 875–895.
- Basilevsky, A. (1994) *Statistical Factor Analysis and Related Methods*, New York: Wiley.
- Benignus, V. A., Muller, K. E., Smith, M. V., Pieper, K. S. and Prah, J. D. (1990) Compensatory tracking in humans with elevated carboxyhemoglobin, *Neurotoxicology and Teratology*, **12**, 105–110.
- Berger, M. P. F. (1986) A Monte Carlo study of the power of alternative tests under the generalized MANOVA model, *Communications in Statistics A*, **15**, 1251–1283.
- Billio, M. and Monfort, A. (1998) Switching state-space models: likelihood function, filtering and smoothing, *Statistical Planning and Inference*, **68**, 65–103.
- Boik, R. J. (1988) The mixed model for multivariate repeated measures: validity conditions and an approximate test, *Psychometrika*, **53**, 469–486.
- Boik, R. J. (1991) Scheffé mixed model for repeated measures a relative efficiency evaluation, *Communications in Statistics - Theory and Methods*, **A20**, 1233–1255.
- Bowden, D. C. (1970) Simultaneous confidence bands for linear regression models, *Journal of the American Statistical Association*, **65**, 413–421.
- Box, G. E. P. (1954a) Some theorems on quadratic forms applied in the study of analysis of variance problems: I. effects of inequality of variance in the one-way classification, *Annals of Mathematical Statistics*, **25**, 290–302.

- Box, G. E. P. (1954b) Some theorems on quadratic forms applied in the study of analysis of variance problems: II. effects of inequality of variance and of correlation between errors in the two-way classification, *Annals of Mathematical Statistics*, **25**, 484–498.
- Box, G. E. P. and Cox, D. R. (1964) An analysis of transformations, *Journal of the Royal Statistical Society B*, **26**, 211–43 (with discussion, 244–252).
- Box, G. E. P. and Cox, D. R. (1984) An analysis of transformations revisited, rebuttal, *Journal of the American Statistical Association*, **77**, 209–210.
- Box, G. E. P., Jenkins, G. M., and Reinsel, G. C. (1994) *Time Series Analysis: Forecasting and Control*, 3rd ed., Englewood Cliffs, NJ: Prentice-Hall.
- Carroll, R. J. and Rupert, D. (1981) On prediction and the power transformation family, *Biometrika*, **68**, 609–615.
- Carroll R. J. and Rupert, D. (1985) Transformations in regression: a robust analysis, *Technometrics*, **27**, 1–12.
- Carroll, R. J. and Rupert, D. (1988) *Transformations and Weighting in Regression*, London: Chapman and Hall.
- Casella, G. (1983) Leverage and regression through the origin, *The American Statistician*, **37**, 147–1512.
- Catellier, D. J. and Muller, K. E. (2000) Tests for Gaussian repeated measures with missing data in small samples, *Statistics in Medicine*, **19**, 1101–1114.
- Catellier, D. J. and Muller, K. E. (2002) Sample size and power considerations in physical activity research. In *Physical Activity Assessments for Health-Related Research*, Chapter 6, G. J. Welk, ed., 93–104. Champaign, IL: Human Kinetics.
- Cochran, W. G. (1934) The distribution of quadratic forms in a normal system with applications to the analysis of variance, *Proceedings of the Cambridge Philosophical Society*, **30**, 178–191.
- Coffey, C. S. and Muller, K. E. (1999) Exact test size and power of a Gaussian error linear model for an internal pilot study, *Statistics in Medicine* **18**, 1199–1214.
- Coffey, C. S. and Muller, K. E. (2000a) Properties of doubly-truncated Gamma variables, *Communications in Statistics A*, **29**, 851–857.
- Coffey, C. S. and Muller, K. E. (2000b) Some distributions and their implications for an internal pilot study with a univariate linear model, *Communications in Statistics A*, **29**, 2677–2691.
- Coffey, C. S. and Muller, K. E. (2001) Controlling test size while gaining the benefits of an internal pilot design, *Biometrics*, **57**, 625–631.
- Coffey, C. S. and Muller, K. E. (2003) Properties of internal pilots with the univariate approach to repeated measures, *Statistics in Medicine*, **22**, 2469–2485.
- Cox, D. R. and Small, N. J. H. (1978) Testing multivariate normality, *Biometrika*, **65**, 263–272.
- Cramer, E. M. and Nicewander, W. A. (1979) Some symmetric, invariant measures of multivariate association, *Psychometrika*, **44**, 43–54.

- Cramér, H. (1946) *Mathematical Methods of Statistics*, Princeton, NJ: University Press.
- Cressie, N. (1991) *Statistics for Spatial Data*, New York: Wiley InterScience.
- D'Agostino, R. B. and Stephens, M. A. (1986) *Goodness-of-Fit-Techniques*, New York: Marcel Dekker.
- Daintith, J. and Nelson, R. D. (1989) *The Penguin Dictionary of Mathematics*, London, UK: Penguin.
- Davies, R. B. (1980) Algorithm AS 155: the distribution of a linear combination of χ^2 random variables, *Applied Statistics*, **29**, 323–333.
- Demidenko, E. (2004) *Mixed Models Theory and Applications*, New York: Wiley.
- Diggle, P. and Kenward, M. G. (1994) Informative drop-out in longitudinal data analysis, *Applied Statistics*, **43**, 49–93.
- Edwards, L. J., Stewart, P. W., Muller, K. E. and Helms, R. W. (2001) Linear equality constraints in the general mixed model, *Biometrics*, **57**, 1185–1190.
- Epps, T. W. (1993) Characteristic functions and their empirical counterparts: geometrical interpretations and applications to statistical inference, *American Statistician*, **47**, 33–38.
- Fairclough, D. L. and Helms, R. H. (1986) A mixed linear model with linear covariance structure: a sensitivity analysis of maximum likelihood estimators, *Journal of Statistical Computation and Simulation*, **25**, 205–236.
- Feller, W. (1968) *An Introduction to Probability Theory and Its Applications*, 3rd ed., Vol. I, New York: Wiley.
- Gallant, A. R. and Fuller, W. A. (1973) Fitting segmented polynomial regression models whose join points have to be estimated, *Journal of the American Statistical Association*, **68**, 144–147.
- Gatsonis, C. and Sampson, A. R. (1989) Multiple correlation: exact power and sample size calculations, *Psychological Bulletin*, **106**, 516–524.
- Geisser, S. and Greenhouse, S. W. (1958) An extension of Box's results on the use of the F distribution in multivariate analysis, *Annals of Mathematical Statistics*, **29**, 885–891.
- Genizi, A. and Soller, M. (1979) Power derivation in an ANOVA model which is intermediate between the “fixed-effects” and the “random-effects” models, *Journal of Statistical Planning and Inference*, **3**, 127–134.
- Gil-Pelaez, J. (1951) Note on the inversion theorem, *Biometrika*, **38**, 481–482.
- Glueck, D. H. and Muller, K. E. (1998) On the trace of a Wishart, *Communications in Statistics A*, **27**, 2137–2141. Corrigenda, 2002, **31**, 159–160.
- Glueck, D. H. and Muller, K. E. (2003) Adjusting power for a baseline covariate in a linear model, *Statistics in Medicine*, **22**, 2535–2551.
- Gnanadesikan, R. and Kettenring, J. R. (1972) Robust estimates, residuals, and outlier detection with multiresponse data, *Biometrics*, **28**, 81–124. Corrigenda, 1972, **28**, 1142.

- Good, I. J. (1963) On the independence of quadratic expressions, *Journal of the Royal Statistical Society B*, **25**, 377–382. Corrigenda, 1966, **28**, 584.
- Good, I. J. (1969) Conditions for a quadratic form to have a chi-squared distribution, *Biometrika*, 1969, **54**, 215–216. Corrigenda, 1970, **57**, 225.
- Goodnight, J. H. (1979) A tutorial on the SWEEP operator, *The American Statistician*, **33**, 149–158.
- Graybill, F. A. (1969) *Introduction to Matrices with Applications in Statistics*. Belmont, CA: Wadsworth.
- Graybill, F. A. (1976) *Theory and Applications of the Linear Model*, North Scituate, MA: Duxbury Press.
- Greenhouse, S. W. and Geisser, S. (1959) On methods in the analysis of profile data, *Psychometrika*, **24**, 95–112.
- Grizzle, J. E. and Allen, D. M. (1969) Analysis of growth and dose response curves, *Biometrics*, **25**, 357–381.
- Gupta, A. K. and Nagar, D. K. (2000) *Matrix Variate Distributions*, Boca Raton, FL: Chapman Hall/CRC.
- Harris, R. J. (1975) *Primer of Multivariate Statistics*, New York: Academic.
- Harvey, J. R. (1972) On expressing moments in terms of cumulants and vice versa, *American Statistician*, **26**, 38–39. Corrigenda, 1973, **27**, 44.
- Harville, D. A. (1974) Bayesian inference for variance components using only error contrasts, *Biometrika*, **61**, 383–385.
- Harville, D. A. (1977) Maximum likelihood approaches to variance component estimation and to related problems, *Journal of the American Statistical Association*, **72**, 320–338. Comments and responses 338–340.
- Harville, D. A. (1990) BLUP (best linear unbiased prediction) and beyond, *Advances in Statistical Methods for Genetic Improvement of Livestock*, pp. 239–276, D. Gianola and K. Hammond Eds., New York: Springer.
- Harville, D. A. (1996) The posterior distribution of the fixed and random effects in a mixed-effects linear model, *Journal of Statistical Computation and Simulation*, **54**, 211–229.
- Helms, R. W. (1988a) Definitions of linear model parameters and hypotheses as functions of $E(Y)$, *Communications in Statistics A*, **17**, 2715–2723.
- Helms, R. W. (1988b) Comparisons of parameter and hypothesis definitions in a general linear model. *Communications in Statistics A*, **17**, 2725–2753.
- Helms, R. W. (1988c) Manipulating statistical hypotheses and tests as Boolean functions, *The American Statistician*, **42**, 253–256.
- Helms, R. W. (1992) Intentionally incomplete longitudinal designs: I. methodology and comparison of some full span designs, *Statistics in Medicine*, **11**, 1889–1913.
- Heitjan, D. F. (1994) Ignorability in general incomplete-data models, *Biometrika*, **81**, 701–708.

- Henderson, C. R. (1953) Estimation of variance and variance components, *Biometrics*, **9**, 226-252.
- Henderson, C. R. (1963) Selection index and expected genetic advance, *Statistical Genetics in Plant Breeding*, National Academy of Sciences, National Research Council Publication #982.
- Hinkley, D. V. and Runger, G. (1984) The analysis of transformed data, *Journal of the American Statistical Association*, **79**, 302-320.
- Hocking, R. R. (1985) *The Analysis of Linear Models*, Monterey, CA: Brooks Cole.
- Hotelling, H. (1951) A generalized t test and measure of multivariate dispersion, *Proceedings of the Second Berkeley Symposium on Mathematical Statistics and Probability*, Berkeley, CA: University of California Press, pp. 23-41.
- Huynh, H. and Feldt, L. S. (1976) Estimation of the Box correction for degrees of freedom from sample data in randomized block and split-plot designs, *Journal of Educational Statistics*, **1**, 69-82.
- Imhof, J. P. (1961) Computing the distribution of quadratic forms in normal variables, *Biometrika*, **48**, 419-426.
- Jackson, G. M. (1991) *A User's Guide to Principal Components*, Chichester, UK: Wiley.
- Jayakar, A. D. (1970) On the detection and estimation of linkage between a locus influencing a quantitative character and a marker locus, *Biometrics*, **26**, 451-464.
- Jennrich, R. I. and Schluchter, M. D. (1985) Unbalanced repeated-measures models with structured covariance matrices, *Biometrics*, **42**, 805-820.
- Jiroutek, M. R., Muller, K. E., Kupper, L. L. and Stewart, P. W. (2003) A new method for choosing sample size for confidence interval-based inferences, *Biometrics*, **59**, 580-590.
- John, S. (1972) The distribution of a statistic used for testing sphericity of normal distributions, *Biometrika*, **59**, 169-173.
- Johnson, N. L. and Kotz, S. (1970) *Continuous Univariate Distributions*, Vol. 2, Boston: Houghton Mifflin.
- Johnson, N. L. and Kotz, S. (1972) *Distributions in Statistics: Continuous Multivariate Distributions*, New York: Wiley.
- Johnson, N. L., Kotz, S. and Balakrishnan, N. (1994) *Continuous Univariate Distributions*, Vol. 1, 2nd ed., New York: Wiley.
- Johnson, N. L., Kotz, S. and Balakrishnan, N. (1995) *Continuous Univariate Distributions*, Vol. 2, 2nd ed., New York: Wiley.
- Johnson, N. L., Kotz, S. and Balakrishnan, N. (1997) *Discrete Multivariate Distributions*, New York: Wiley.
- Johnson, N. L., Kotz, S. and Balakrishnan, N. (1995) *Continuous Univariate Distributions*, Vol. 1, 2nd ed., New York: Wiley.
- Jolliffe, I. T. (2002) *Principal Components Analysis*, 2nd ed., New York: Springer-Verlag.

- Jöreskog, K. G. (1993) *Applied Factor Analysis in the Natural Sciences*, Cambridge: Cambridge University Press.
- Kackar, R. N. and Harville, D. A. (1984) Approximations for standard errors of estimators of fixed and random effects in mixed linear models, *Journal of the American Statistical Association*, **79**, 853–862.
- Kendall, M. G. and Stuart, A. (1977) *The Advanced Theory of Statistics*, Vol. 1, 4th ed., London: Charles Griffin and Co.
- Kenward, M. G. and Roger, J. H. (1997) Small sample inference for fixed effects from restricted maximum likelihood, *Biometrics*, **53**, 983–997.
- Khatri, C. G. (1966) A note on a MANOVA model applied to problems in growth curves, *Annals of the Institute of Statistical Mathematics*, **18**, 75–86.
- Khuri, A. I., Mathew, T. and Sinha B. K. (1998) *Statistical Tests for Mixed Linear Models*, New York: Wiley.
- Kim, H., Gribbin, M. J., Muller, K. E. and Taylor, D. J. (2006) Analytic, computational and approximate forms for ratios of noncentral and central Gaussian quadratic forms, *Journal of Computational and Graphical Statistics*, **15**, 443–459.
- Kirk, R. E. (1995) *Experimental Design: Procedures for the Behavioral Sciences*, 3rd ed., Monterey, CA: Brooks Cole.
- Kleinbaum, D. G., Kupper, L. L., Muller, K. E., and Nizam, A. (1998) *Applied Regression Analysis and Other Multivariable Methods*, 3rd ed., Boston: Duxbury Press.
- Kotz, S., Balakrishnan, N. and Johnson, N. L., (2000) *Continuous Multivariate Distributions*, Vol. 1, 2nd ed., New York: Wiley.
- Kotz, S. and Nadarajah, S. (2004) *Multivariate t Distributions and Their Applications*, Cambridge: Cambridge University Press.
- Kshirsagar, A. M., and Smith, W. B. (1995) *Growth Curves*, New York: Marcel Dekker.
- Kuhfeld, W. F. (1986) A note on Roy's largest root, *Psychometrika*, **51**, 479–481.
- Kvålseth, T. O. (1985) Cautionary note about R^2 , *The American Statistician*, **39**, 279–285.
- Laird, N. M. and Ware, J. (1982) Random-effects models for longitudinal data, *Biometrics*, **38**, 963–974.
- Lancaster, P. (1969) *Theory of Matrices*, New York: Academic Press.
- Lawley, D. N. (1938) A generalization of Fisher's z test, *Biometrika*, **30**, 180–187.
- Lenth, R. V. (2001) Some practical guidelines for effective sample size determination, *American Statistician*, **55**, 187–193.
- Lindgren, B. W. (1976) *Statistical Theory*, 3rd ed., New York: Macmillan.
- Lindstrom, M. J. and Bates, D. F. (1988) Newton-Raphson and EM algorithms for linear mixed effects models for repeated measures data, *Journal of the American Statistical Association*, **83**, 1014–1022. Corrigendum, 1994, **89**, 1572.

- Lipsitz, S. R., Garrett I., Fitzmaurice, M., Ibrahim, J. G., Gelber, R. and Lipshultz, S. (2002) Parameter estimation in longitudinal studies with outcome-dependent follow-up, *Biometrics*, **58**, 621–630.
- Littel, R. (2003) Analysis of unbalanced mixed models: a case study comparison of ANOVA vs REML/GLS, *Journal of Agricultural, Biological, and Environmental Statistics*, **7**, 472–490.
- Little, R. J. A. (1992), Regression with missing X's: a review, *Journal of the American Statistical Association*, **87**, 1227–1237.
- Little, R. J. A. and Rubin, D. B. (2002) *Statistical Analysis with Missing Data*, 2nd ed., New York: Wiley.
- Loperfido, N. (1999) A conditional normality test based on the conditional probability integral transformation, *Statistica Applicata*, **11**, 281–291.
- Loynes, R. M. (1966) On idempotent matrices, *Annals of Mathematical Statistics*, **37**, 491–494.
- Lukacs, E. (1983) *Developments in Characteristic Function Theory*, London: Charles Griffin and Co.
- Mecklin, C. J. and Mundfrom, D. J. (2005) A Monte Carlo comparison of the type I and type II error rates of test of multivariate normality, *Journal of Statistical Computation and Simulation*, **75**, 93–107.
- Magnus, J. R. (1978) Maximum likelihood estimation of the GLS model with unknown parameters in the disturbance covariance matrix, *Journal of Econometrics*, **7**, 281–312.
- Magnus, J. R. and Neudecker, H. (1988) *Matrix Differential Calculus With Applications in Statistics and Econometrics*, New York: Wiley.
- Malkovich, J. F. and Afifi, A. A. (1973) On tests for multivariate normality, *Journal of the American Statistical Association*, **68**, 176–179.
- Mardia, K. V. (1970) Measures of multivariate skewness and kurtosis with applications, *Biometrika*, **57**, 519–530.
- Mardia, K. V. (1975) [Algorithm AS 84] Measures of multivariate skewness and kurtosis, *Applied Statistics*, **24**, 262–265.
- Mardia, K. V., Kent, J. T. and Bibby, J. M. (1979) *Multivariate Analysis*, London: Academic Press.
- McCullagh, P. and Nelder, J. A. (1989) *Generalized Linear Models*, 2nd ed., London: Chapman & Hall.
- McDonald, R. P. (1985) *Factor Analysis and Related Methods*, NJ: Lawrence Erlbaum.
- McKeon, J. J. (1974) F approximations to the distribution of Hotelling's T_0^2 , *Biometrika*, **61**, 381–383.
- Mecklin, C. J. and Mundfrom, D. J. (2005) A Monte Carlo comparison of the type I and type II error rates of tests of multivariate normality, *Journal of Statistical Computation and Simulation*, **75**, 93-107.

- Meng, X. and Rubin, D. B. (1991) Using EM to obtain asymptotic variance-covariance matrices: the SEM algorithm, *Journal of the American Statistical Association*, **86**, 899–909.
- Mitra, S. K. (1970) A density-free approach to the matrix variate beta distribution, *Sankhya A*, **32**, 81–88.
- Morrison, D. F. (1990) *Multivariate Statistical Methods*, 3rd ed., New York: McGraw-Hill.
- Muirhead, R. J. (1984) *Aspects of Multivariate Statistical Theory*, New York: Wiley.
- Muller, K. E. (1982) Understanding canonical correlation through the general linear model and principal components, *American Statistician*, **36**, 342–354.
- Muller, K. E. (1998) A new F approximation for the Pillai-Bartlett trace under H_0 , *Journal of Computational and Graphical Statistics*, **7**, 131–137.
- Muller, K. E. and Barton, C. N. (1989) Approximate power for repeated-measures ANOVA lacking sphericity, *Journal of the American Statistical Association*, **84**, 549–555. Corrigenda, 1991, **86**, 255–256.
- Muller, K. E., Barton, C. N., and Benignus, V.A. (1984) Recommendations for appropriate statistical practice in toxicologic experiments, *Neurotoxicology*, **5**, 113–126.
- Muller, K. E. and Benignus, V. A. (1992) Increasing scientific power with statistical power, *Neurotoxicology and Teratology*, **14**, 211–219.
- Muller, K. E. and Chi, Y. (2006) Extending the central Wishart to high dimension, low sample size, with implications for simulations, *manuscript in review*.
- Muller, K. E., Edwards, L. J., Simpson, S. and Taylor D. T. (2006) Accurate test size and power in small samples of repeated measures without missing data, *manuscript in review*.
- Muller, K. E. and Fetterman, B. A. (2002) *Regression and ANOVA: An Integrated Approach Using SAS® Software*, Cary, NC: SAS Institute.
- Muller, K. E., LaVange, L. M., Ramey, S. L. and Ramey, C. T. (1992) Power calculations for general linear multivariate models including repeated measures applications, *Journal of the American Statistical Association*, **87**, 1209–1226.
- Muller, K. E., and Pasour, V. B. (1997) Bias in linear model power and sample size due to estimating variance, *Communications in Statistics A*, **26**, 839–851.
- Muller, K. E. and Peterson, B. L. (1984) Practical methods for computing power in testing the multivariate general linear hypothesis, *Computational Statistics and Data Analysis*, **2**, 143–158.
- O'Brien, R. G. (1979) A general ANOVA method for robust tests of additive models for variances, *Journal of the American Statistical Association*, **74**, 877–880.
- O'Brien, R. G. and Muller, K. E. (1993) A unified approach to statistical power for t -tests to multivariate models. In *Applied Analysis of Variance in Behavioral Sciences*, Chapter 8, pp. 297–344, L. K. Edwards Ed., New York: Marcel Dekker.

- Odell, P. L. and Feiveson, A. H. (1966) A numerical procedure to generate a sample covariance matrix, *Journal of the American Statistical Association*, **61**, 199–203. Corrigendum, 1249–1250, with acknowledgment of priority to Wijsman, 1959.
- Olson, C. L. (1974) Comparative robustness of six tests in multivariate analysis of variance, *Journal of the American Statistical Association*, **69**, 894–908.
- Olson, C. L. (1976) On choosing a test statistic in multivariate analysis, *Psychological Bulletin*, **83**, 579–586.
- Olson, C. L. (1979) Practical considerations in choosing a MANOVA test statistic: a rejoinder to Stevens, *Psychological Bulletin*, **86**, 1350–1352.
- Pan, W. K. (2003) *Multilevel Spatial and Statistical Analyses to Examine the Relationship between Population and Environment: A Case Study of the Ecuadorian Amazon*, Dissertation, Department of Biostatistics, University of North Carolina, Chapel Hill.
- Park, T., Park, J. K. and Davis, C. S. (2001) Effects of covariance model assumptions on hypothesis tests for repeated measurements: analysis of ovarian hormone data and pituitary-pteryomaxillary distance data, *Statistics in Medicine*, **20**, 2441–2453.
- Patterson, H. D. and Thompson, R. (1971) Recovery of inter-block information when block sizes are unequal, *Biometrika*, **58**, 545–554.
- Pillai, K. C. S. (1955) Some new test criteria in multivariate analysis, *Annals of Mathematical Statistics*, **26**, 117–121.
- Pillai, K. C. S. (1956) On the distribution of the largest or the smallest root of a matrix in multivariate analysis, *Biometrika*, **43**, 122–127.
- Pizer, S. M., Chen, J. Z., Fletcher, P. T., Friedman, Y., Fritsch, D. S., Gash, A. G., Glotzer, J. M., Jiroutek, M. R., Joshi, S., Lu, C., Muller, K. E., Thall, A., Tracton, G., Yushkevich, P. and Chaney, E. L. (2003) Deformable m-reps for 3D medical image segmentation, *International Journal of Computer Vision*, **55**, 85–106.
- Potthoff, R. F. and Roy, S. N. (1964) A generalized multivariate analysis of variance model useful especially for growth curve problems, *Biometrics*, **51**, 313–326.
- Rao, C. R. (1951) An asymptotic expansion of the distribution of Wilks' criterion, *Bulletin of the Institute of International Statistics*, **XXXIII**, 177–180.
- Rao, C. R. (1965) The theory of least squares when the parameters are stochastic and its application to the analysis of growth curves, *Biometrika*, **52**, 447–468.
- Rao, C. R. (1973) *Linear Statistical Inference and Its Applications*, 2nd ed., New York: Wiley.
- Raudenbush, S. W. and Bryk, A. S. (2002) *Hierarchical Linear Models: Applications and Data Analysis Methods*, Newbury Park, CA: Sage Publications.
- Rocke, D. M. (1989) Bootstrap Bartlett adjustment in seemingly unrelated regression, *Journal of the American Statistical Association*, **84**, 598–601.
- Rousseeuw, P. J. and Van Zomeren, B. C. (1990) Unmasking multivariate outliers and leverage points, *Journal of the American Statistical Association*, **85**, 633–639. Also see comments and rejoinder, 640–651.

- Roy, S. N. (1957) *Some Aspects of Multivariate Analysis*, New York: Wiley.
- Rubin, D. R. (1976) Inference and missing data, *Biometrika*, **63**, 581-590.
- Sampson, A. R. (1974) A tale of two regressions, *Journal of the American Statistical Association*, **69**, 682-689.
- SAS Institute (1999) *SAS/IML User's Guide, Version 8*, Cary, NC: SAS Institute, Inc.
- Satterthwaite, F. E. (1946) An approximate distribution of estimates of variance components, *Biometrics Bulletin*, **2**, 110-114.
- Schaalje, G. B., McBride J. B. and Fellingham, G. W. (2003) Adequacy of approximations to distributions of test statistics in complex mixed linear models, *Journal of Agricultural, Biological, and Environmental Statistics*, **7**, 512-524.
- Schott, J. R. (2005) *Matrix Analysis for Statistics*, 2nd ed. New York: Wiley.
- Scott, A. and Wild, C. (1991) Transformations and R^2 , *The American Statistician*, **45**, 127-129.
- Searle, S. R. (1971) *Linear Models*, New York: Wiley.
- Searle, S. R. (1982) *Matrix Algebra Useful for Statistics*, New York: Wiley.
- Searle, S. R., Casella, G. and McCulloch, C. M. (1992) *Variance Components*, New York: Wiley.
- Shanbhag, D. N. (1966) On the independence of quadratic forms, *Journal of the Royal Statistical Society B*, **28**, 582-583.
- Siegel, A. F. (1979) The noncentral chi-squared distribution with zero degrees of freedom and testing for uniformity, *Biometrika*, **66**, 381-386.
- Small, N. J. H. (1980) Marginal skewness and kurtosis in testing multivariate normality, *Applied Statistics*, **29**, 85-87.
- Small, N. J. H. (1988) In *The Encyclopedia of Statistical Sciences*, Vol. 9, pp. 95-100, N. L. Johnson and S. Kotz, Eds., New York: Wiley.
- Smith, P. L. (1979) Splines as a useful and convenient statistical tool, *The American Statistician*, **33**, 57-62.
- Soller, M. and Genizi, A. (1978) The efficiency of experimental designs for the detection of linkage between a marker locus and a locus affecting a quantitative trait in segregating populations, *Biometrics*, **34**, 37-55.
- Srivastava, V. K., and Giles, D. E. A. (1987) *Seemingly Unrelated Regression Equation Models: Estimation and Inference*, New York: Marcel Dekker.
- Stevens, S. S. (1946) On the theory of scales of measurement, *Science*, **103**, 677-680.
- Stevens, S. S. (1951) Mathematics, measurement, and psychophysics. In *Handbook of Experimental Psychology*, S. S. Stevens, ed., New York: Wiley.
- Stewart, P. W. (1987) Line segment confidence limits, *Biometrics*, **43**, 629-640.
- Stewart, P. W. (1991) The graphical advantages of finite interval confidence band procedures, *Communications in Statistics A*, **20**, 3975-3994.
- Stewart, P. W. (2000) Estimation and comparison of growth and dose-response curves in the presence of purposeful censoring. In *Handbook of Statistics, Vol. 18:*

- Bioenvironmental and Public Health Statistics*, P. K. Sen, P. K. and C. R. Rao, Eds. Ch.10, pp. 325–354, Amsterdam and New York: North-Holland/Elsevier.
- Stroup, W. W. (2003) Power analysis based on spatial effects mixed models: a tool for comparing design and analysis strategies in the presence of spatial variability, *Journal of Agricultural, Biological, and Environmental Statistics*, **7**, 491–511.
- Taylor, D. J. and Muller, K. E. (1995) Computing confidence bounds for power and sample size of the general linear univariate model, *American Statistician*, **49**, 43–47.
- Taylor, D. J. and Muller, K. E. (1996) Bias in linear model power and sample size calculations due to estimating noncentrality, *Communications in Statistics A*, **25**, 1595–1610.
- Thode, H. C. Jr. (2002) *Testing for Normality*, New York: Marcel Dekker.
- Timm, N. H. (1975) *Multivariate Analysis with Applications in Education and Psychology*, Monterey, CA: Brooks/Cole.
- Timm, N. H. (2002) *Applied Multivariate Analysis*, New York: Springer-Verlag.
- Velleman, P. F. and Wilkinson, L. (1993) Nominal, ordinal, interval, and ratio typologies are misleading, *American Statistician*, **47**, 65–72.
- Verbeke, G. and Molenberghs, G. (2000) *Linear Mixed Models for Longitudinal Data*, New York: Springer.
- Vonesh, E. F. and Chinchilli, V. M. (1997) *Linear and Nonlinear Models for the Analysis of Repeated Measurements*, New York: Marcel Dekker.
- Wackerly, D. D., Mendenhall, W., and Scheaffer, R. L. (1996) *Mathematical Statistics with Applications*, Belmont, CA: Duxbury Press, 5th ed.
- Weisstein, E. W. (2003) *CRC Concise Encyclopedia of Mathematics*, 2nd ed., Boca Raton, FL: Chapman & Hall/CRC.
- Widaman, K. F. (1993) Common factor analysis versus principal component analysis: differential bias in representing model parameters?, *Multivariate Behavioral Research*, **28**, 263–311.
- Widaman, K. F. (2004) Common factors vs. components: principals and principles, errors and misconceptions, presentation at *Factor Analysis at 100: Historical Developments and Future Directions*, conference at the University of North Carolina, Chapel Hill, May 2004.
- Wijisman, R. A. (1959) Applications of a certain representation of the Wishart matrix, *Annals of Mathematical Statistics*, **30**, 597–601.
- Willet, J. B. and Singer, J. D. (1988) Another cautionary note about R^2 : its use in weighted least-squares regression analysis, *The American Statistician*, **42**, 236–238.
- Wishart, J. (1928) The generalized product moment distribution in samples from a normal multivariate population, *Biometrika*, **20A**, 32–52.
- Wittes, J. and Brittain, E. (1990) The role of internal pilot studies in increasing the efficiency of clinical trials, *Statistics in Medicine*, **9**, 65–72.

- Wolfinger, R., Tobias, R. and Sall, J. (1994) Computing Gaussian likelihoods and their derivatives for general linear mixed models, *SIAM Journal on Scientific and Statistical Computing*, **15**, 1294–1310.
- Zucker, D. M. and Denne, J. (2002) Sample size redetermination for repeated measures studies, *Biometrics*, **58**, 548–559.
- Zucker, D. M., Lieberman, O. and Manor, O. (2000) Improved small sample inference in the mixed linear model: Bartlett correction and adjusted likelihood, *Journal of the Royal Statistical Society B*, **62**, 827–838.

Index

A

- Acceptance region, 51
- Added last tests, 356–360
- Added-in-order tests, 356–360
- Algebraic multiplicity, 20
- Approximate weighted least squares (AWLS), 82
- Associative laws for matrix algebra, 9
- Autoregressive covariance matrix, 38

B

- Balanced design, 105
- Balanced random coefficient model, 283
- Basis, 11
- Best linear unbiased estimator (BLUE), 222, 227
- Best linear unbiased predictor (BLUP), 287
- Bickel-Doksum transformation, 128
- Bilinear form, 10
- Biorthogonal matrices, 10
- Block diagonal matrix, 3
- Boolean algebra, 44
- Bonferroni correction, 67, 77, 107–108
- Box-Cox power transformation, 128
- Breakdown point, 162

C

- Canonical correlation, 64
- Casewise missing, 74–75
- Cauchy-Schwartz inequality, 135
- Characteristic function (CF)
 - Vector, 119
 - Matrix, 138
- Chi-square distribution, 169–170
- Cholesky factor matrix, 25, 72
- Coefficient of determination, 47
- Column, extracting from matrix, 4
- Commensurate, 101
- Common factor model, 332
- Commutative laws for matrix algebra, 9
- Complete design, 105
- Composite acceptance region, 301
- Composite alternative, 301
- Component hypothesis, 301
- Composite hypothesis, 300
- Compound symmetric covariance matrix, 37–38
- Concatenating matrices, 4
- Conditional distribution, 132
- Confidence
 - Band, finite interval, 340
 - Coefficient, 52
 - Intervals, 52
 - Regions, 52, 73, 306, 332
- Conform, conformation of matrices, 5
- Congruent matrices, 21
- Consistent system of equations, 17

Constituent matrix decomposition, 23
 Continuous
 Data, 103
 Random variable, including
 absolutely, 117
 Correlation matrix, 33, 135
 Covariance design matrix, 92
 Covariance (dispersion) matrix, 8, 25,
 33, 134
 sample value, 34
 For Gaussian data, also see *Wishart*
 Critical region (rejection region), 51
 Critical value, 51
 Cross product, 6
 Cumulant, generating function (CGF),
 124, 125
 Cumulative distribution function (CDF),
 116

D

Data matrix, 34
 Determinant, 14–15
 Diagonal of matrix, 2
 Diagonal matrix, 2, 14, 22
 Direct product (matrix multiplication), 7
 Direct-product (matrix)
 Gaussian distribution, 156, 157, 349
 Direct sum, 7
 Discrete Fourier transform, 273
 Discrete random variable, 117
 Distributions, 116, 117
 Distributive laws for matrix algebra, 9
 Dot product, 6
 Doubly Multivariate outcomes, 101

E

Eigenanalysis, 18, 19, 20, 21
 Eigenvalue, 19
 Eigenvector, 20
 Error sum of squares, 49, 65

Essence matrix, 218
 Estimable, 212, 249
 Exact weighted least squares, 238
 Exchangeable observations, 105
 Excludes an intercept, 46
 Expected value, 132

F

Factor matrix, 24
 Factor analysis, 332
 Fisher Scoring algorithm, 284
 Full rank, 12–15

G

Gaussian (nonsingular or singular,
 Standard or not, central or noncentral)
 Scalar, 142
 Vector (multivariate), 143–144
 General matrix, 156
 Direct-product matrix, 156
 General linear (null) hypothesis, 60
 General linear mixed model, 92
 With Gaussian errors, 96
 General linear multivariate model, 56
 With Gaussian errors, 59
 Generalized inverse matrix, 14, 16–17,
 23, 25
 Geometric multiplicity, 20
 Generalized general linear model
 (GGLM), 80–81
 linearly equivalent, 238
 General linear model (GLM), univariate
 Analysis of Variance (ANOVA)
 models, 42
 Definition, 40
 Coefficient of determination, 47, 48
 Concepts, 39
 Estimated errors, 65
 Estimation, 39
 Estimators, 209, 246

- Explicit and implicit restrictions, 40, 231
 - Full rank, 40, 56
 - Inference, 39
 - Least squares assumptions, 41
 - Less than full rank, 40, 224
 - Noncentrality parameter, 44
 - Numerical methods, 39
 - One-to-one linear transformations, 267
 - Predicted values, 65
 - Primary parameters, 40, 209
 - Secondary parameters, 43, 209
 - With Gaussian errors, 41
 - General linear model (GLM),
 - multivariate
 - Estimators, 246
 - Primary parameters, 59, 246
 - Secondary parameters, 59, 246
 - Growth curve model (GCM), 84–86, 89, 266–277
- H*
- Hat matrix, 49, 65
 - Horizontal direct product, 6
 - Horizontally concatenated, 4
 - Hotelling-Lawley trace statistic, 61, 67, 319
 - Hypothesis
 - A priori parameter, 289
 - Full and constrained model, 45
 - General linear hypothesis (GLH), 43–44, 46, 48, 59, 289
 - Post hoc parameters, 289
 - Sum of squares, 49, 65
 - Testability, 290–295
 - Tests, multivariate, 66–73, Chapter 16
 - Tests, univariate, 50–51, Chapter 15
- I*
- Idempotent, 10
 - Identity matrix, 3
 - Ignorable (missing data), 76
 - Independence
 - Mutual (total), 130
 - Observation, 101
 - Pairwise, 130
 - Statistical, 130
 - Independent sampling unit (ISU), 101
 - Inner product, 6, 10
 - Interval scale, 103
 - Inverse (matrix operator), 15
 - Inverse Wishart, 205
 - Iterated approximate weighted least squares (ITAWLS), 82, 263
- J*
- Joint cumulative distribution function, 117
- K*
- Knots, see *Spline*
 - Kurtosis, 164
 - Kronecker covariance, 245
 - Kronecker design, 245
 - Kronecker (direct) product, 7
- L*
- Less than full rank, 12, 15, 48, 95
 - Least squares estimator, 224
 - Linear covariance structure, 282
 - Linear mixed models tests, 341–343
 - Linear transformation, including full rank (nonsingular) and less than full rank (singular), 126
 - Linearly equivalent, 210, 246, 269
 - Linearly dependent, 11
 - Linearly independent, 11
 - Linear (scale) invariance, 305, 317
 - Location invariance, 305, 317
 - Lower triangular matrix, 3

- M*
- Mahalanobis distance, 62, 162
 - Marginal distribution, 128–129
 - Masking, 162
 - Matrix addition and subtraction, 5
 - Matrix, defined, 1
 - Mean
 - Population, 32, 132
 - Sample, 34
 - Also see *moments*
 - Missing
 - Data, 75
 - At random (MAR), 76
 - Completely at random (MCAR), 75
 - Model notation, summary, 111
 - Moment generating function (MGF)
 - Vector, 123
 - Matrix, 138
 - Moments, 133
 - Moore-Penrose generalized inverse, 16
 - Multiplication
 - Direct product of two matrices, 7
 - Elementwise for two matrices, 5
 - Matrix with matrix, 2, 6
 - Horizontal direct product, 6
 - Scalar with matrix, 5
 - Multivariate analysis of variance (MANOVA), 107
 - Multivariate association, 70
 - Multivariate outcomes, 101
 - Multivariate quadratic form, 193
 - Multivariate tests, including multivariate approach to repeated measures (MULTIREP), 59, 61, 70, 318
- N*
- Negative definite, 22
 - Negative semidefinite, 22
 - Nominal scale, 103
 - Noncentral chi-square distribution, 170
 - Noncentral Wishart distribution, 350
 - Noncentrality parameter, 44, 59
 - Nonnegative definite, 22
 - Nonpositive definite, 22
 - Nonsingular matrix, 15
 - Nonsingular Gaussian vector, 144, 146
 - Nonsymmetric square matrices, 12
 - Nonzero eigenvalues, 21
 - Null hypothesis, 44
- O*
- Observational unit, 101
 - Operations in matrix algebra, rules for, 8–9
 - Ordinal scale, 103
 - Ordinary least squares (OLS), 263
 - Orthogonal matrix, 10
 - Orthonormal matrix, 10
 - Outer product, 6, 10
- P*
- Partitioned matrix, 3, 4, 30, 33
 - Permutation matrix, 30
 - Pillai-Bartlett trace, ANOVA analog statistic, 67, 319
 - Polynomial growth curve, 269
 - Positive definite, 22
 - Positive semidefinite, 22
 - Power of a test, 51
 - Power function, 51
 - Predicted values, 49, 65
 - Principal component, 36
 - Projection, projection matrix, 27–28
 - Probability density function (PDF), 116, 118, 129
 - Projection matrix, 28
- Q*
- Quadratic form
 - Matrix expression, 10
 - Random, 174

univariate Gaussian, 174
 multivariate Gaussian, 193

R

Random deviation, 94
 Random vector, 143
 Rank of a matrix, defined, 12
 Rank of a vector space, 12
 Ratio scale, 103
 Rejection region (critical region), 51
 Repeated measures, 101
 Residuals, 49, 65
 Restricted maximum likelihood (REML), 285
 Restricted linear model (explicit or implicit), 231
 Row, extracting from matrix, 4
 Roy's largest root, the union-intersection principle statistic, 66, 319
 Robust distance, 162

S

Sample mean (vector), 34
 Sample SSCP matrix, 34
 Sample covariance matrix, 34
 Scalar multiplication of a matrix, 5
 Scalar, defined, 1
 Scale (linear) invariance, 305, 317
 Secondary parameter, 43, 209, 246
 Semidefinite, 22
 Seemingly unrelated regressions, 83
 Shift parameter, 44, 59
 Similar matrices, 21
 Simple matrix, 20
 Singular matrix, 15
 Singular value decomposition, 13, 26
 Singular Gaussian vector, 144
 Size of a test, 51
 Skewness, 164
 Spans an intercept, 46

Spectral decomposition, 13, 36
 Spectrum (of a matrix), 19
 Sphericity, spherical distribution, 37
 Spline, 235
 Square matrix, 2
 Squared error loss function, 223
 Stacked by column, 4
 Standard Gaussian: see Gaussian
 Stieltjes integral, 119
 Subpopulation mean, 94
 Sums of squares and cross products (SSCP) matrix, 34
 Sylvester's Law of Inertia, 21
 Symmetric matrix, 2

T

Testable, 290, 312
 Testability of hypothesis, 312, 316
 Test statistic, 51
 Trace, matrix, 5
 Transpose, 2
 Translation, 126
 Triangular matrix (upper, lower), 3
 Type I error, 51
 Type II error, 51

U

Uncorrelated, 136
 Uniformly minimum variance estimator, (UMVUE), 222, 227
 Union-intersection test, 301
 Univariate approach to repeated measures (UNIREP) tests, 61, 67
 Unstructured covariance matrix, 37
 Upper triangular matrix, 3

V

Variance, 134
 Vector space, defined, 1, 11

Vertically concatenated, 4

W

Weighted least squares (WLS), 82

Wilks lambda, likelihood test, 67, 319

Wishart distribution, 193, 350

Z

Zero matrix, 3