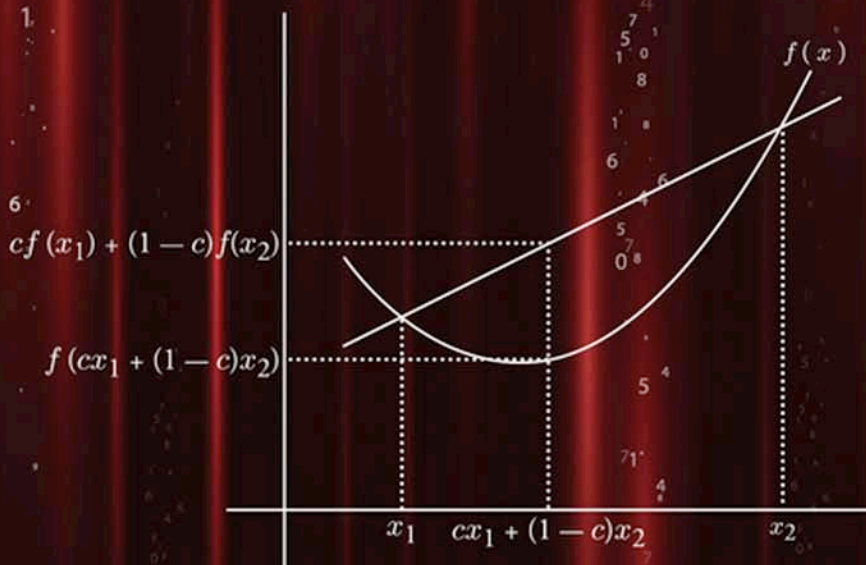


Wiley Series in Probability and Statistics

THIRD EDITION

MATRIX ANALYSIS FOR STATISTICS

James R. Schott



WILEY

MATRIX ANALYSIS FOR STATISTICS

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: David J. Balding, Noel A. C. Cressie, Garrett M. Fitzmaurice,
Geof H. Givens, Harvey Goldstein, Geert Molenberghs, David W. Scott,
Adrian F. M. Smith, Ruey S. Tsay, Sanford Weisberg

Editors Emeriti: J. Stuart Hunter, Iain M. Johnstone, Joseph B. Kadane,
Jozef L. Teugels

A complete list of the titles in this series appears at the end of this volume.

MATRIX ANALYSIS FOR STATISTICS

Third Edition

JAMES R. SCHOTT

WILEY

Copyright © 2017 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/permissions>.

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 formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Names: Schott, James R., 1955- author.

Title: Matrix analysis for statistics / James R. Schott.

Description: Third edition. | Hoboken, New Jersey : John Wiley & Sons, 2016.

| Includes bibliographical references and index.

Identifiers: LCCN 2016000005 | ISBN 9781119092483 (cloth) | ISBN 9781119092469 (epub)

Subjects: LCSH: Matrices. | Mathematical statistics.

Classification: LCC QA188 .S24 2016 | DDC 512.9/434--dc23 LC record available at <http://lccn.loc.gov/2016000005>

Cover image courtesy of GettyImages/Alexmumu.

Typeset in 10/12pt TimesLTStd by SPi Global, Chennai, India

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

To Susan, Adam, and Sarah

CONTENTS

Preface	xi
About the Companion Website	xv
1 A Review of Elementary Matrix Algebra	1
1.1 Introduction, 1	
1.2 Definitions and Notation, 1	
1.3 Matrix Addition and Multiplication, 2	
1.4 The Transpose, 3	
1.5 The Trace, 4	
1.6 The Determinant, 5	
1.7 The Inverse, 9	
1.8 Partitioned Matrices, 12	
1.9 The Rank of a Matrix, 14	
1.10 Orthogonal Matrices, 15	
1.11 Quadratic Forms, 16	
1.12 Complex Matrices, 18	
1.13 Random Vectors and Some Related Statistical Concepts, 19	
Problems, 29	
2 Vector Spaces	35
2.1 Introduction, 35	
2.2 Definitions, 35	
2.3 Linear Independence and Dependence, 42	

2.4	Matrix Rank and Linear Independence, 45	
2.5	Bases and Dimension, 49	
2.6	Orthonormal Bases and Projections, 53	
2.7	Projection Matrices, 58	
2.8	Linear Transformations and Systems of Linear Equations, 65	
2.9	The Intersection and Sum of Vector Spaces, 73	
2.10	Oblique Projections, 76	
2.11	Convex Sets, 80	
	Problems, 85	
3	Eigenvalues and Eigenvectors	95
3.1	Introduction, 95	
3.2	Eigenvalues, Eigenvectors, and Eigenspaces, 95	
3.3	Some Basic Properties of Eigenvalues and Eigenvectors, 99	
3.4	Symmetric Matrices, 106	
3.5	Continuity of Eigenvalues and Eigenprojections, 114	
3.6	Extremal Properties of Eigenvalues, 116	
3.7	Additional Results Concerning Eigenvalues Of Symmetric Matrices, 123	
3.8	Nonnegative Definite Matrices, 129	
3.9	Antieigenvalues and Antieigenvectors, 141	
	Problems, 144	
4	Matrix Factorizations and Matrix Norms	155
4.1	Introduction, 155	
4.2	The Singular Value Decomposition, 155	
4.3	The Spectral Decomposition of a Symmetric Matrix, 162	
4.4	The Diagonalization of a Square Matrix, 169	
4.5	The Jordan Decomposition, 173	
4.6	The Schur Decomposition, 175	
4.7	The Simultaneous Diagonalization of Two Symmetric Matrices, 178	
4.8	Matrix Norms, 184	
	Problems, 191	
5	Generalized Inverses	201
5.1	Introduction, 201	
5.2	The Moore–Penrose Generalized Inverse, 202	
5.3	Some Basic Properties of the Moore–Penrose Inverse, 205	
5.4	The Moore–Penrose Inverse of a Matrix Product, 211	
5.5	The Moore–Penrose Inverse of Partitioned Matrices, 215	
5.6	The Moore–Penrose Inverse of a Sum, 219	
5.7	The Continuity of the Moore–Penrose Inverse, 222	
5.8	Some Other Generalized Inverses, 224	

5.9	Computing Generalized Inverses, 232 Problems, 238	
6	Systems of Linear Equations	247
6.1	Introduction, 247	
6.2	Consistency of a System of Equations, 247	
6.3	Solutions to a Consistent System of Equations, 251	
6.4	Homogeneous Systems of Equations, 258	
6.5	Least Squares Solutions to a System of Linear Equations, 260	
6.6	Least Squares Estimation For Less Than Full Rank Models, 266	
6.7	Systems of Linear Equations and The Singular Value Decomposition, 271	
6.8	Sparse Linear Systems of Equations, 273 Problems, 278	
7	Partitioned Matrices	285
7.1	Introduction, 285	
7.2	The Inverse, 285	
7.3	The Determinant, 288	
7.4	Rank, 296	
7.5	Generalized Inverses, 298	
7.6	Eigenvalues, 302 Problems, 307	
8	Special Matrices and Matrix Operations	315
8.1	Introduction, 315	
8.2	The Kronecker Product, 315	
8.3	The Direct Sum, 323	
8.4	The Vec Operator, 323	
8.5	The Hadamard Product, 329	
8.6	The Commutation Matrix, 339	
8.7	Some Other Matrices Associated With the Vec Operator, 346	
8.8	Nonnegative Matrices, 351	
8.9	Circulant and Toeplitz Matrices, 363	
8.10	Hadamard and Vandermonde Matrices, 369 Problems, 373	
9	Matrix Derivatives and Related Topics	387
9.1	Introduction, 387	
9.2	Multivariable Differential Calculus, 387	
9.3	Vector and Matrix Functions, 390	
9.4	Some Useful Matrix Derivatives, 396	
9.5	Derivatives of Functions of Patterned Matrices, 400	

9.6	The Perturbation Method, 402	
9.7	Maxima and Minima, 409	
9.8	Convex and Concave Functions, 413	
9.9	The Method of Lagrange Multipliers, 417	
	Problems, 423	
10	Inequalities	433
10.1	Introduction, 433	
10.2	Majorization, 433	
10.3	Cauchy-Schwarz Inequalities, 444	
10.4	Hölder's Inequality, 446	
10.5	Minkowski's Inequality, 450	
10.6	The Arithmetic-Geometric Mean Inequality, 452	
	Problems, 453	
11	Some Special Topics Related to Quadratic Forms	457
11.1	Introduction, 457	
11.2	Some Results on Idempotent Matrices, 457	
11.3	Cochran's Theorem, 462	
11.4	Distribution of Quadratic Forms in Normal Variates, 465	
11.5	Independence of Quadratic Forms, 471	
11.6	Expected Values of Quadratic Forms, 477	
11.7	The Wishart Distribution, 485	
	Problems, 496	
	References	507
	Index	513

PREFACE

As the field of statistics has developed over the years, the role of matrix methods has evolved from a tool through which statistical problems could be more conveniently expressed to an absolutely essential part in the development, understanding, and use of the more complicated statistical analyses that have appeared in recent years. As such, a background in matrix analysis has become a vital part of a graduate education in statistics. Too often, the statistics graduate student gets his or her matrix background in bits and pieces through various courses on topics such as regression analysis, multivariate analysis, linear models, stochastic processes, and so on. An alternative to this fragmented approach is an entire course devoted to matrix methods useful in statistics. This text has been written with such a course in mind. It also could be used as a text for an advanced undergraduate course with an unusually bright group of students and should prove to be useful as a reference for both applied and research statisticians.

Students beginning in a graduate program in statistics often have their previous degrees in other fields, such as mathematics, and so initially their statistical backgrounds may not be all that extensive. With this in mind, I have tried to make the statistical topics presented as examples in this text as self-contained as possible. This has been accomplished by including a section in the first chapter which covers some basic statistical concepts and by having most of the statistical examples deal with applications which are fairly simple to understand; for instance, many of these examples involve least squares regression or applications that utilize the simple concepts of mean vectors and covariance matrices. Thus, an introductory statistics course should provide the reader of this text with a sufficient background in statistics. An additional prerequisite is an undergraduate course in matrices or linear algebra, while a calculus background is necessary for some portions of the book, most notably, Chapter 8.

By selectively omitting some sections, all nine chapters of this book can be covered in a one-semester course. For instance, in a course targeted at students who end their educational careers with the masters degree, I typically omit Sections 2.10, 3.5, 3.7, 4.8, 5.4-5.7, and 8.6, along with a few other sections.

Anyone writing a book on a subject for which other texts have already been written stands to benefit from these earlier works, and that certainly has been the case here. The texts by Basilevsky (1983), Graybill (1983), Healy (1986), and Searle (1982), all books on matrices for statistics, have helped me, in varying degrees, to formulate my ideas on matrices. Graybill's book has been particularly influential, since this is the book that I referred to extensively, first as a graduate student, and then in the early stages of my research career. Other texts which have proven to be quite helpful are Horn and Johnson (1985, 1991), Magnus and Neudecker (1988), particularly in the writing of Chapter 8, and Magnus (1988).

I wish to thank several anonymous reviewers who offered many very helpful suggestions, and Mark Johnson for his support and encouragement throughout this project. I am also grateful to the numerous students who have alerted me to various mistakes and typos in earlier versions of this book. In spite of their help and my diligent efforts at proofreading, undoubtedly some mistakes remain, and I would appreciate being informed of any that are spotted.

JIM SCHOTT

Orlando, Florida

PREFACE TO THE SECOND EDITION

The most notable change in the second edition is the addition of a chapter on results regarding matrices partitioned into a 2×2 form. This new chapter, which is Chapter 7, has the material on the determinant and inverse that was previously given as a section in Chapter 7 of the first edition. Along with the results on the determinant and inverse of a partitioned matrix, I have added new material in this chapter on the rank, generalized inverses, and eigenvalues of partitioned matrices.

The coverage of eigenvalues in Chapter 3 has also been expanded. Some additional results such as Weyl's Theorem have been included, and in so doing, the last section of Chapter 3 of the first edition has now been replaced by two sections.

Other smaller additions, including both theorems and examples, have been made elsewhere throughout the book. Over 100 new exercises have been added to the problems sets.

The writing of a second edition of this book has also given me the opportunity to correct mistakes in the first edition. I would like to thank those readers who have

pointed out some of these errors as well as those that have offered suggestions for improvement to the text.

JIM SCHOTT

*Orlando, Florida
September 2004*

PREFACE TO THE THIRD EDITION

The third edition of this text maintains the same organization that was present in the previous editions. The major changes involve the addition of new material. This includes the following additions.

1. A new chapter, now Chapter 10, on inequalities has been added. Numerous inequalities such as Cauchy-Schwarz, Hadamard, and Jensen's, already appear in the earlier editions, but there are many important ones that are missing, and some of these are given in the new chapter. Highlighting this chapter is a fairly substantial section on majorization and some of the inequalities that can be developed from this concept.
2. A new section on oblique projections has been added to Chapter 2. The previous editions only covered orthogonal projections.
3. A new section on antieigenvalues and antieigenvectors has been added to Chapter 3.

Numerous other smaller additions have been made throughout the text. These include some additional theorems, the proofs of some results that previously had been given without proof, and some more examples involving statistical applications. Finally, more than 70 new problems have been added to the end-of-chapter problem sets.

JIM SCHOTT

*Orlando, Florida
December 2015*

ABOUT THE COMPANION WEBSITE

This book is accompanied by a companion website:

www.wiley.com/go/Schott/MatrixAnalysis3e

The instructor's website includes:

- A solutions manual with solutions to selected problems

The student's website includes:

- A solutions manual with odd-numbered solutions to selected problems



1

A REVIEW OF ELEMENTARY MATRIX ALGEBRA

1.1 INTRODUCTION

In this chapter, we review some of the basic operations and fundamental properties involved in matrix algebra. In most cases, properties will be stated without proof, but in some cases, when instructive, proofs will be presented. We end the chapter with a brief discussion of random variables and random vectors, expected values of random variables, and some important distributions encountered elsewhere in the book.

1.2 DEFINITIONS AND NOTATION

Except when stated otherwise, a scalar such as α will represent a real number. A matrix A of size $m \times n$ is the $m \times n$ rectangular array of scalars given by

$$A = \begin{bmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{bmatrix},$$

and sometimes it is simply identified as $A = (a_{ij})$. Sometimes it also will be convenient to refer to the (i, j) th element of A , as $(A)_{ij}$; that is, $a_{ij} = (A)_{ij}$. If $m = n$,

then A is called a square matrix of order m , whereas A is referred to as a rectangular matrix when $m \neq n$. An $m \times 1$ matrix

$$\mathbf{a} = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \end{bmatrix}$$

is called a column vector or simply a vector. The element a_i is referred to as the i th component of \mathbf{a} . A $1 \times n$ matrix is called a row vector. The i th row and j th column of the matrix A will be denoted by $(A)_i$ and $(A)_{.j}$, respectively. We will usually use capital letters to represent matrices and lowercase bold letters for vectors.

The diagonal elements of the $m \times m$ matrix A are $a_{11}, a_{22}, \dots, a_{mm}$. If all other elements of A are equal to 0, A is called a diagonal matrix and can be identified as $A = \text{diag}(a_{11}, \dots, a_{mm})$. If, in addition, $a_{ii} = 1$ for $i = 1, \dots, m$ so that $A = \text{diag}(1, \dots, 1)$, then the matrix A is called the identity matrix of order m and will be written as $A = I_m$ or simply $A = I$ if the order is obvious. If $A = \text{diag}(a_{11}, \dots, a_{mm})$ and b is a scalar, then we will use A^b to denote the diagonal matrix $\text{diag}(a_{11}^b, \dots, a_{mm}^b)$. For any $m \times m$ matrix A , D_A will denote the diagonal matrix with diagonal elements equal to those of A , and for any $m \times 1$ vector \mathbf{a} , $D_{\mathbf{a}}$ denotes the diagonal matrix with diagonal elements equal to the components of \mathbf{a} ; that is, $D_A = \text{diag}(a_{11}, \dots, a_{mm})$ and $D_{\mathbf{a}} = \text{diag}(a_1, \dots, a_m)$.

A triangular matrix is a square matrix that is either an upper triangular matrix or a lower triangular matrix. An upper triangular matrix is one that has all of its elements below the diagonal equal to 0, whereas a lower triangular matrix has all of its elements above the diagonal equal to 0. A strictly upper triangular matrix is an upper triangular matrix that has each of its diagonal elements equal to 0. A strictly lower triangular matrix is defined similarly.

The i th column of the $m \times m$ identity matrix will be denoted by \mathbf{e}_i ; that is, \mathbf{e}_i is the $m \times 1$ vector that has its i th component equal to 1 and all of its other components equal to 0. When the value of m is not obvious, we will make it more explicit by writing \mathbf{e}_i as $\mathbf{e}_{i,m}$. The $m \times m$ matrix whose only nonzero element is a 1 in the (i, j) th position will be identified as E_{ij} .

The scalar zero is written 0, whereas a vector of zeros, called a null vector, will be denoted by $\mathbf{0}$, and a matrix of zeros, called a null matrix, will be denoted by (0) . The $m \times 1$ vector having each component equal to 1 will be denoted by $\mathbf{1}_m$ or simply $\mathbf{1}$ when the size of the vector is obvious.

1.3 MATRIX ADDITION AND MULTIPLICATION

The sum of two matrices A and B is defined if they have the same number of rows and the same number of columns; in this case,

$$A + B = (a_{ij} + b_{ij}).$$

The product of a scalar α and a matrix A is

$$\alpha A = A\alpha = (\alpha a_{ij}).$$

The premultiplication of the matrix B by the matrix A is defined only if the number of columns of A equals the number of rows of B . Thus, if A is $m \times p$ and B is $p \times n$, then $C = AB$ will be the $m \times n$ matrix which has its (i, j) th element, c_{ij} , given by

$$c_{ij} = (A)_{i.}(B)_{.j} = \sum_{k=1}^p a_{ik}b_{kj}.$$

A similar definition exists for BA , the postmultiplication of B by A , if the number of columns of B equals the number of rows of A . When both products are defined, we will not have, in general, $AB = BA$. If the matrix A is square, then the product AA , or simply A^2 , is defined. In this case, if we have $A^2 = A$, then A is said to be an idempotent matrix.

The following basic properties of matrix addition and multiplication in Theorem 1.1 are easy to verify.

Theorem 1.1 Let α and β be scalars and A , B , and C be matrices. Then, when the operations involved are defined, the following properties hold:

- (a) $A + B = B + A$.
- (b) $(A + B) + C = A + (B + C)$.
- (c) $\alpha(A + B) = \alpha A + \alpha B$.
- (d) $(\alpha + \beta)A = \alpha A + \beta A$.
- (e) $A - A = A + (-A) = (0)$.
- (f) $A(B + C) = AB + AC$.
- (g) $(A + B)C = AC + BC$.
- (h) $(AB)C = A(BC)$.

1.4 THE TRANSPOSE

The transpose of an $m \times n$ matrix A is the $n \times m$ matrix A' obtained by interchanging the rows and columns of A . Thus, the (i, j) th element of A' is a_{ji} . If A is $m \times p$ and B is $p \times n$, then the (i, j) th element of $(AB)'$ can be expressed as

$$\begin{aligned} ((AB)')_{ij} &= (AB)_{ji} = (A)_{j.}(B)_{.i} = \sum_{k=1}^p a_{jk}b_{ki} \\ &= (B')_i.(A')_{.j} = (B'A')_{ij}. \end{aligned}$$

Thus, evidently $(AB)' = B'A'$. This property along with some other results involving the transpose are summarized in Theorem 1.2.

Theorem 1.2 Let α and β be scalars and A and B be matrices. Then, when defined, the following properties hold:

- (a) $(\alpha A)' = \alpha A'$.
- (b) $(A')' = A$.
- (c) $(\alpha A + \beta B)' = \alpha A' + \beta B'$.
- (d) $(AB)' = B'A'$.

If A is $m \times m$, that is, A is a square matrix, then A' is also $m \times m$. In this case, if $A = A'$, then A is called a symmetric matrix, whereas A is called a skew-symmetric if $A = -A'$.

The transpose of a column vector is a row vector, and in some situations, we may write a matrix as a column vector times a row vector. For instance, the matrix E_{ij} defined in Section 1.2 can be expressed as $E_{ij} = e_i e_j'$. More generally, $e_{i,m} e_{j,n}'$ yields an $m \times n$ matrix having 1, as its only nonzero element, in the (i, j) th position, and if A is an $m \times n$ matrix, then

$$A = \sum_{i=1}^m \sum_{j=1}^n a_{ij} e_{i,m} e_{j,n}'.$$

1.5 THE TRACE

The trace is a function that is defined only on square matrices. If A is an $m \times m$ matrix, then the trace of A , denoted by $\text{tr}(A)$, is defined to be the sum of the diagonal elements of A ; that is,

$$\text{tr}(A) = \sum_{i=1}^m a_{ii}.$$

Now if A is $m \times n$ and B is $n \times m$, then AB is $m \times m$ and

$$\begin{aligned} \text{tr}(AB) &= \sum_{i=1}^m (AB)_{ii} = \sum_{i=1}^m (A)_{i \cdot} (B)_{\cdot i} = \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ji} \\ &= \sum_{j=1}^n \sum_{i=1}^m b_{ji} a_{ij} = \sum_{j=1}^n (B)_{j \cdot} (A)_{\cdot j} \\ &= \sum_{j=1}^n (BA)_{jj} = \text{tr}(BA). \end{aligned}$$

This property of the trace, along with some others, is summarized in Theorem 1.3.

Theorem 1.3 Let α be a scalar and A and B be matrices. Then, when the appropriate operations are defined, we have the following properties:

- (a) $\text{tr}(A') = \text{tr}(A)$.
- (b) $\text{tr}(\alpha A) = \alpha \text{tr}(A)$.
- (c) $\text{tr}(A + B) = \text{tr}(A) + \text{tr}(B)$.
- (d) $\text{tr}(AB) = \text{tr}(BA)$.
- (e) $\text{tr}(A'A) = 0$ if and only if $A = (0)$.

1.6 THE DETERMINANT

The determinant is another function defined on square matrices. If A is an $m \times m$ matrix, then its determinant, denoted by $|A|$, is given by

$$\begin{aligned} |A| &= \sum (-1)^{f(i_1, \dots, i_m)} a_{1i_1} a_{2i_2} \cdots a_{mi_m} \\ &= \sum (-1)^{f(i_1, \dots, i_m)} a_{i_1 1} a_{i_2 2} \cdots a_{i_m m}, \end{aligned}$$

where the summation is taken over all permutations (i_1, \dots, i_m) of the set of integers $(1, \dots, m)$, and the function $f(i_1, \dots, i_m)$ equals the number of transpositions necessary to change (i_1, \dots, i_m) to an increasing sequence of components, that is, to $(1, \dots, m)$. A transposition is the interchange of two of the integers. Although f is not unique, it is uniquely even or odd, so that $|A|$ is uniquely defined. Note that the determinant produces all products of m terms of the elements of the matrix A such that exactly one element is selected from each row and each column of A .

Using the formula for the determinant, we find that $|A| = a_{11}$ when $m = 1$. If A is 2×2 , we have

$$|A| = a_{11}a_{22} - a_{12}a_{21},$$

and when A is 3×3 , we get

$$\begin{aligned} |A| &= a_{11}a_{22}a_{33} + a_{12}a_{23}a_{31} + a_{13}a_{21}a_{32} \\ &\quad - a_{11}a_{23}a_{32} - a_{12}a_{21}a_{33} - a_{13}a_{22}a_{31}. \end{aligned}$$

The following properties of the determinant in Theorem 1.4 are fairly straightforward to verify using the definition of a determinant.

Theorem 1.4 If α is a scalar and A is an $m \times m$ matrix, then the following properties hold:

- (a) $|A'| = |A|$.
- (b) $|\alpha A| = \alpha^m |A|$.

- (c) If A is a diagonal matrix, then $|A| = a_{11} \cdots a_{mm} = \prod_{i=1}^m a_{ii}$.
- (d) If all elements of a row (or column) of A are zero, $|A| = 0$.
- (e) The interchange of two rows (or columns) of A changes the sign of $|A|$.
- (f) If all elements of a row (or column) of A are multiplied by α , then the determinant is multiplied by α .
- (g) The determinant of A is unchanged when a multiple of one row (or column) is added to another row (or column).
- (h) If two rows (or columns) of A are proportional to one another, $|A| = 0$.

An alternative expression for $|A|$ can be given in terms of the cofactors of A . The minor of the element a_{ij} , denoted by m_{ij} , is the determinant of the $(m-1) \times (m-1)$ matrix obtained after removing the i th row and j th column from A . The corresponding cofactor of a_{ij} , denoted by A_{ij} , is then given as $A_{ij} = (-1)^{i+j} m_{ij}$.

Theorem 1.5 For any $i = 1, \dots, m$, the determinant of the $m \times m$ matrix A can be obtained by expanding along the i th row,

$$|A| = \sum_{j=1}^m a_{ij} A_{ij}, \quad (1.1)$$

or expanding along the i th column,

$$|A| = \sum_{j=1}^m a_{ji} A_{ji}. \quad (1.2)$$

Proof. We will just prove (1.1), as (1.2) can easily be obtained by applying (1.1) to A' . We first consider the result when $i = 1$. Clearly

$$\begin{aligned} |A| &= \sum (-1)^{f(i_1, \dots, i_m)} a_{1i_1} a_{2i_2} \cdots a_{mi_m} \\ &= a_{11} b_{11} + \cdots + a_{1m} b_{1m}, \end{aligned}$$

where

$$a_{1j} b_{1j} = \sum (-1)^{f(i_1, \dots, i_m)} a_{1i_1} a_{2i_2} \cdots a_{mi_m},$$

and the summation is over all permutations for which $i_1 = j$. Since $(-1)^{f(j, i_2, \dots, i_m)} = (-1)^{j-1} (-1)^{f(i_2, \dots, i_m)}$, this implies that

$$b_{1j} = \sum (-1)^{j-1} (-1)^{f(i_2, \dots, i_m)} a_{2i_2} \cdots a_{mi_m},$$

where the summation is over all permutations (i_2, \dots, i_m) of $(1, \dots, j-1, j+1, \dots, m)$. If C is the $(m-1) \times (m-1)$ matrix obtained from A by deleting its 1st row and j th column, then b_{1j} can be written

$$\begin{aligned}
b_{1j} &= (-1)^{j-1} \sum (-1)^{f(i_1, \dots, i_{m-1})} c_{1i_1} \cdots c_{m-1i_{m-1}} = (-1)^{j-1} |C| \\
&= (-1)^{j-1} m_{1j} = (-1)^{1+j} m_{1j} = A_{1j},
\end{aligned}$$

where the summation is over all permutations (i_1, \dots, i_{m-1}) of $(1, \dots, m-1)$ and m_{1j} is the minor of a_{1j} . Thus,

$$|A| = \sum_{j=1}^m a_{1j} b_{1j} = \sum_{j=1}^m a_{1j} A_{1j},$$

as is required. To prove (1.1) when $i > 1$, let D be the $m \times m$ matrix for which $(D)_{1\cdot} = (A)_{i\cdot}$, $(D)_j = (A)_{j-1\cdot}$, for $j = 2, \dots, i$, and $(D)_j = (A)_j$ for $j = i+1, \dots, m$. Then $A_{ij} = (-1)^{i-1} D_{1j}$, $a_{ij} = d_{1j}$ and $|A| = (-1)^{i-1} |D|$. Thus, since we have already established (1.1) when $i = 1$, we have

$$|A| = (-1)^{i-1} |D| = (-1)^{i-1} \sum_{j=1}^m d_{1j} D_{1j} = \sum_{j=1}^m a_{ij} A_{ij},$$

and so the proof is complete. \square

Our next result indicates that if the cofactors of a row or column are matched with the elements from a different row or column, the expansion reduces to 0.

Theorem 1.6 If A is an $m \times m$ matrix and $k \neq i$, then

$$\sum_{j=1}^m a_{ij} A_{kj} = \sum_{j=1}^m a_{ji} A_{jk} = 0. \quad (1.3)$$

Example 1.1 We will find the determinant of the 5×5 matrix given by

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

Using the cofactor expansion formula on the first column of A , we obtain

$$|A| = 2 \begin{vmatrix} 0 & 3 & 0 & 0 \\ 0 & 2 & 2 & 0 \\ 0 & 1 & 1 & 1 \\ 1 & 2 & 2 & 1 \end{vmatrix},$$

and then using the same expansion formula on the first column of this 4×4 matrix, we get

$$|A| = 2(-1) \begin{vmatrix} 3 & 0 & 0 \\ 2 & 2 & 0 \\ 1 & 1 & 1 \end{vmatrix}.$$

Because the determinant of the 3×3 matrix above is 6, we have

$$|A| = 2(-1)(6) = -12.$$

Consider the $m \times m$ matrix C whose columns are given by the vectors $\mathbf{c}_1, \dots, \mathbf{c}_m$; that is, we can write $C = (\mathbf{c}_1, \dots, \mathbf{c}_m)$. Suppose that, for some $m \times 1$ vector $\mathbf{b} = (b_1, \dots, b_m)'$ and $m \times m$ matrix $A = (\mathbf{a}_1, \dots, \mathbf{a}_m)$, we have

$$\mathbf{c}_1 = A\mathbf{b} = \sum_{i=1}^m b_i \mathbf{a}_i.$$

Then, if we find the determinant of C by expanding along the first column of C , we get

$$\begin{aligned} |C| &= \sum_{j=1}^m c_{j1} C_{j1} = \sum_{j=1}^m \left(\sum_{i=1}^m b_i a_{ji} \right) C_{j1} \\ &= \sum_{i=1}^m b_i \left(\sum_{j=1}^m a_{ji} C_{j1} \right) = \sum_{i=1}^m b_i |(\mathbf{a}_i, \mathbf{c}_2, \dots, \mathbf{c}_m)|, \end{aligned}$$

so that the determinant of C is a linear combination of m determinants. If B is an $m \times m$ matrix and we now define $C = AB$, then by applying the previous derivation on each column of C , we find that

$$\begin{aligned} |C| &= \left| \left(\sum_{i_1=1}^m b_{i_1 1} \mathbf{a}_{i_1}, \dots, \sum_{i_m=1}^m b_{i_m m} \mathbf{a}_{i_m} \right) \right| \\ &= \sum_{i_1=1}^m \cdots \sum_{i_m=1}^m b_{i_1 1} \cdots b_{i_m m} |(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m})| \\ &= \sum b_{i_1 1} \cdots b_{i_m m} |(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m})|, \end{aligned}$$

where this final sum is only over all permutations of $(1, \dots, m)$, because Theorem 1.4(h) implies that

$$|(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m})| = 0$$

if $i_j = i_k$ for any $j \neq k$. Finally, reordering the columns in $|(\mathbf{a}_{i_1}, \dots, \mathbf{a}_{i_m})|$ and using Theorem 1.4(e), we have

$$|C| = \sum b_{i_1 1} \cdots b_{i_m m} (-1)^{f(i_1, \dots, i_m)} |(\mathbf{a}_1, \dots, \mathbf{a}_m)| = |B||A|.$$

This very useful result is summarized in Theorem 1.7.

Theorem 1.7 If both A and B are square matrices of the same order, then

$$|AB| = |A||B|.$$

1.7 THE INVERSE

An $m \times m$ matrix A is said to be a nonsingular matrix if $|A| \neq 0$ and a singular matrix if $|A| = 0$. If A is nonsingular, a nonsingular matrix denoted by A^{-1} and called the inverse of A exists, such that

$$AA^{-1} = A^{-1}A = I_m. \quad (1.4)$$

This inverse is unique because, if B is another $m \times m$ matrix satisfying the inverse formula (1.4) for A , then $BA = I_m$, and so

$$B = BI_m = BAA^{-1} = I_m A^{-1} = A^{-1}.$$

The following basic properties of the matrix inverse in Theorem 1.8 can be easily verified by using (1.4).

Theorem 1.8 If α is a nonzero scalar, and A and B are nonsingular $m \times m$ matrices, then the following properties hold:

- (a) $(\alpha A)^{-1} = \alpha^{-1} A^{-1}$.
- (b) $(A')^{-1} = (A^{-1})'$.
- (c) $(A^{-1})^{-1} = A$.
- (d) $|A^{-1}| = |A|^{-1}$.
- (e) If $A = \text{diag}(a_{11}, \dots, a_{mm})$, then $A^{-1} = \text{diag}(a_{11}^{-1}, \dots, a_{mm}^{-1})$.
- (f) If $A = A'$, then $A^{-1} = (A^{-1})'$.
- (g) $(AB)^{-1} = B^{-1}A^{-1}$.

As with the determinant of A , the inverse of A can be expressed in terms of the cofactors of A . Let $A_{\#}$, called the adjoint of A , be the transpose of the matrix of cofactors of A ; that is, the (i, j) th element of $A_{\#}$ is A_{ji} , the cofactor of a_{ji} . Then

$$AA_{\#} = A_{\#}A = \text{diag}(|A|, \dots, |A|) = |A|I_m,$$

because $(A)_{i \cdot} (A_{\#})_{\cdot i} = (A_{\#})_{i \cdot} (A)_{\cdot i} = |A|$ follows directly from (1.1) and (1.2), and $(A)_{i \cdot} (A_{\#})_{\cdot j} = (A_{\#})_{i \cdot} (A)_{\cdot j} = 0$, for $i \neq j$ follows from (1.3). The equation above then yields the relationship

$$A^{-1} = |A|^{-1} A_{\#}$$

when $|A| \neq 0$. Thus, for instance, if A is a 2×2 nonsingular matrix, then

$$A^{-1} = |A|^{-1} \begin{bmatrix} a_{22} & -a_{12} \\ -a_{21} & a_{11} \end{bmatrix}.$$

Similarly when $m = 3$, we get $A^{-1} = |A|^{-1}A_{\#}$, where

$$A_{\#} = \begin{bmatrix} a_{22}a_{33} - a_{23}a_{32} & -(a_{12}a_{33} - a_{13}a_{32}) & a_{12}a_{23} - a_{13}a_{22} \\ -(a_{21}a_{33} - a_{23}a_{31}) & a_{11}a_{33} - a_{13}a_{31} & -(a_{11}a_{23} - a_{13}a_{21}) \\ a_{21}a_{32} - a_{22}a_{31} & -(a_{11}a_{32} - a_{12}a_{31}) & a_{11}a_{22} - a_{12}a_{21} \end{bmatrix}.$$

The relationship between the inverse of a matrix product and the product of the inverses, given in Theorem 1.8(g), is a very useful property. Unfortunately, such a nice relationship does not exist between the inverse of a sum and the sum of the inverses. We do, however, have Theorem 1.9 which is sometimes useful.

Theorem 1.9 Suppose A and B are nonsingular matrices, with A being $m \times m$ and B being $n \times n$. For any $m \times n$ matrix C and any $n \times m$ matrix D , it follows that if $A + CBD$ is nonsingular, then

$$(A + CBD)^{-1} = A^{-1} - A^{-1}C(B^{-1} + DA^{-1}C)^{-1}DA^{-1}.$$

Proof. The proof simply involves verifying that $(A + CBD)(A + CBD)^{-1} = I_m$ for $(A + CBD)^{-1}$ given above. We have

$$\begin{aligned} & (A + CBD)\{A^{-1} - A^{-1}C(B^{-1} + DA^{-1}C)^{-1}DA^{-1}\} \\ &= I_m - C(B^{-1} + DA^{-1}C)^{-1}DA^{-1} + CBDA^{-1} \\ &\quad - CBDA^{-1}C(B^{-1} + DA^{-1}C)^{-1}DA^{-1} \\ &= I_m - C\{(B^{-1} + DA^{-1}C)^{-1} - B \\ &\quad + BDA^{-1}C(B^{-1} + DA^{-1}C)^{-1}\}DA^{-1} \\ &= I_m - C\{B(B^{-1} + DA^{-1}C)(B^{-1} + DA^{-1}C)^{-1} - B\}DA^{-1} \\ &= I_m - C\{B - B\}DA^{-1} = I_m, \end{aligned}$$

and so the result follows. \square

The expression given for $(A + CBD)^{-1}$ in Theorem 1.9 involves the inverse of the matrix $B^{-1} + DA^{-1}C$. It can be shown (see Problem 7.12) that the conditions of the theorem guarantee that this inverse exists. If $m = n$ and C and D are identity matrices, then we obtain Corollary 1.9.1 of Theorem 1.9.

Corollary 1.9.1 Suppose that A, B and $A + B$ are all $m \times m$ nonsingular matrices. Then

$$(A + B)^{-1} = A^{-1} - A^{-1}(B^{-1} + A^{-1})^{-1}A^{-1}.$$

We obtain Corollary 1.9.2 of Theorem 1.9 when $n = 1$.

Corollary 1.9.2 Let A be an $m \times m$ nonsingular matrix. If \mathbf{c} and \mathbf{d} are both $m \times 1$ vectors and $A + \mathbf{cd}'$ is nonsingular, then

$$(A + \mathbf{cd}')^{-1} = A^{-1} - A^{-1}\mathbf{cd}'A^{-1}/(1 + \mathbf{d}'A^{-1}\mathbf{c}).$$

Example 1.2 Theorem 1.9 can be particularly useful when m is larger than n and the inverse of A is fairly easy to compute. For instance, suppose we have $A = I_5$,

$$B = \begin{bmatrix} 1 & 1 \\ 1 & 2 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 0 \\ 2 & 1 \\ -1 & 1 \\ 0 & 2 \\ 1 & 1 \end{bmatrix}, \quad D' = \begin{bmatrix} 1 & -1 \\ -1 & 2 \\ 0 & 1 \\ 1 & 0 \\ -1 & 1 \end{bmatrix},$$

from which we obtain

$$G = A + CBD = \begin{bmatrix} 1 & 1 & 1 & 1 & 0 \\ -1 & 6 & 4 & 3 & 1 \\ -1 & 2 & 2 & 0 & 1 \\ -2 & 6 & 4 & 3 & 2 \\ -1 & 4 & 3 & 2 & 2 \end{bmatrix}.$$

It is somewhat tedious to compute the inverse of this 5×5 matrix directly. However, the calculations in Theorem 1.9 are fairly straightforward. Clearly, $A^{-1} = I_5$ and

$$B^{-1} = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix},$$

so that

$$(B^{-1} + DA^{-1}C) = \begin{bmatrix} 2 & -1 \\ -1 & 1 \end{bmatrix} + \begin{bmatrix} -2 & 0 \\ 3 & 4 \end{bmatrix} = \begin{bmatrix} 0 & -1 \\ 2 & 5 \end{bmatrix}$$

and

$$(B^{-1} + DA^{-1}C)^{-1} = \begin{bmatrix} 2.5 & 0.5 \\ -1 & 0 \end{bmatrix}.$$

Thus, we find that

$$G^{-1} = I_5 - C(B^{-1} + DA^{-1}C)^{-1}D$$

$$= \begin{bmatrix} -1 & 1.5 & -0.5 & -2.5 & 2 \\ -3 & 3 & -1 & -4 & 3 \\ 3 & -2.5 & 1.5 & 3.5 & -3 \\ 2 & -2 & 0 & 3 & -2 \\ -1 & 0.5 & -0.5 & -1.5 & 2 \end{bmatrix}.$$

1.8 PARTITIONED MATRICES

Occasionally we will find it useful to partition a given matrix into submatrices. For instance, suppose A is $m \times n$ and the positive integers m_1, m_2, n_1, n_2 are such that $m = m_1 + m_2$ and $n = n_1 + n_2$. Then one way of writing A as a partitioned matrix is

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

where A_{11} is $m_1 \times n_1$, A_{12} is $m_1 \times n_2$, A_{21} is $m_2 \times n_1$, and A_{22} is $m_2 \times n_2$. That is, A_{11} is the matrix consisting of the first m_1 rows and n_1 columns of A , A_{12} is the matrix consisting of the first m_1 rows and last n_2 columns of A , and so on. Matrix operations can be expressed in terms of the submatrices of the partitioned matrix. For example, suppose B is an $n \times p$ matrix partitioned as

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix},$$

where B_{11} is $n_1 \times p_1$, B_{12} is $n_1 \times p_2$, B_{21} is $n_2 \times p_1$, B_{22} is $n_2 \times p_2$, and $p = p_1 + p_2$. Then the premultiplication of B by A can be expressed in partitioned form as

$$AB = \begin{bmatrix} A_{11}B_{11} + A_{12}B_{21} & A_{11}B_{12} + A_{12}B_{22} \\ A_{21}B_{11} + A_{22}B_{21} & A_{21}B_{12} + A_{22}B_{22} \end{bmatrix}.$$

Matrices can be partitioned into submatrices in other ways besides this 2×2 partitioned form. For instance, we could partition only the columns of A , yielding the expression

$$A = [A_1 \ A_2],$$

where A_1 is $m \times n_1$ and A_2 is $m \times n_2$. A more general situation is one in which the rows of A are partitioned into r groups and the columns of A are partitioned into c groups so that A can be written as

$$A = \begin{bmatrix} A_{11} & A_{12} & \cdots & A_{1c} \\ A_{21} & A_{22} & \cdots & A_{2c} \\ \vdots & \vdots & & \vdots \\ A_{r1} & A_{r2} & \cdots & A_{rc} \end{bmatrix},$$

where the submatrix A_{ij} is $m_i \times n_j$ and the integers m_1, \dots, m_r and n_1, \dots, n_c are such that

$$\sum_{i=1}^r m_i = m \quad \text{and} \quad \sum_{j=1}^c n_j = n.$$

This matrix A is said to be in block diagonal form if $r = c$, A_{ii} is a square matrix for each i , and A_{ij} is a null matrix for all i and j for which $i \neq j$. In this case, we will write $A = \text{diag}(A_{11}, \dots, A_{rr})$; that is,

$$\text{diag}(A_{11}, \dots, A_{rr}) = \begin{bmatrix} A_{11} & (0) & \cdots & (0) \\ (0) & A_{22} & \cdots & (0) \\ \vdots & \vdots & & \vdots \\ (0) & (0) & \cdots & A_{rr} \end{bmatrix}.$$

Example 1.3 Suppose we wish to compute the transpose product AA' , where the 5×5 matrix A is given by

$$A = \begin{bmatrix} 1 & 0 & 0 & 1 & 1 \\ 0 & 1 & 0 & 1 & 1 \\ 0 & 0 & 1 & 1 & 1 \\ -1 & -1 & -1 & 2 & 0 \\ -1 & -1 & -1 & 0 & 2 \end{bmatrix}.$$

The computation can be simplified by observing that A may be written as

$$A = \begin{bmatrix} I_3 & \mathbf{1}_3 \mathbf{1}_2' \\ -\mathbf{1}_2 \mathbf{1}_3' & 2I_2 \end{bmatrix}.$$

As a result, we have

$$\begin{aligned} AA' &= \begin{bmatrix} I_3 & \mathbf{1}_3 \mathbf{1}_2' \\ -\mathbf{1}_2 \mathbf{1}_3' & 2I_2 \end{bmatrix} \begin{bmatrix} I_3 & -\mathbf{1}_3 \mathbf{1}_2' \\ \mathbf{1}_2 \mathbf{1}_3' & 2I_2 \end{bmatrix} \\ &= \begin{bmatrix} I_3 + \mathbf{1}_3 \mathbf{1}_2' \mathbf{1}_2 \mathbf{1}_3' & -\mathbf{1}_3 \mathbf{1}_2' + 2\mathbf{1}_3 \mathbf{1}_2' \\ -\mathbf{1}_2 \mathbf{1}_3' + 2\mathbf{1}_2 \mathbf{1}_3' & \mathbf{1}_2 \mathbf{1}_3' \mathbf{1}_3 \mathbf{1}_2' + 4I_2 \end{bmatrix} \\ &= \begin{bmatrix} I_3 + 2\mathbf{1}_3 \mathbf{1}_3' & \mathbf{1}_3 \mathbf{1}_2' \\ \mathbf{1}_2 \mathbf{1}_3' & 3\mathbf{1}_2 \mathbf{1}_2' + 4I_2 \end{bmatrix} \\ &= \begin{bmatrix} 3 & 2 & 2 & 1 & 1 \\ 2 & 3 & 2 & 1 & 1 \\ 2 & 2 & 3 & 1 & 1 \\ 1 & 1 & 1 & 7 & 3 \\ 1 & 1 & 1 & 3 & 7 \end{bmatrix}. \end{aligned}$$

1.9 THE RANK OF A MATRIX

Our initial definition of the rank of an $m \times n$ matrix A is given in terms of submatrices. We will see an alternative equivalent definition in terms of the concept of linearly independent vectors in Chapter 2. Most of the material we include in this section can be found in more detail in texts on elementary linear algebra such as Andrilli and Hecker (2010) and Poole (2015).

In general, any matrix formed by deleting rows or columns of A is called a submatrix of A . The determinant of an $r \times r$ submatrix of A is called a minor of order r . For instance, for an $m \times m$ matrix A , we have previously defined what we called the minor of a_{ij} ; this is an example of a minor of order $m - 1$. Now the rank of a nonnull $m \times n$ matrix A is r , written $\text{rank}(A) = r$, if at least one of its minors of order r is nonzero while all minors of order $r + 1$ (if there are any) are zero. If A is a null matrix, then $\text{rank}(A) = 0$. If $\text{rank}(A) = \min(m, n)$, then A is said to have full rank. In particular, if $\text{rank}(A) = m$, A has full row rank, and if $\text{rank}(A) = n$, A has full column rank.

The rank of a matrix A is unchanged by any of the following operations, called elementary transformations:

- (a) The interchange of two rows (or columns) of A .
- (b) The multiplication of a row (or column) of A by a nonzero scalar.
- (c) The addition of a scalar multiple of a row (or column) of A to another row (or column) of A .

Thus, the definition of the rank of A is sometimes given as the number of nonzero rows in the reduced row echelon form of A .

Any elementary transformation of A can be expressed as the multiplication of A by a matrix referred to as an elementary transformation matrix. An elementary transformation of the rows of A will be given by the premultiplication of A by an elementary transformation matrix, whereas an elementary transformation of the columns corresponds to a postmultiplication. Elementary transformation matrices are nonsingular, and any nonsingular matrix can be expressed as the product of elementary transformation matrices. Consequently, we have Theorem 1.10.

Theorem 1.10 Let A be an $m \times n$ matrix, B be an $m \times m$ matrix, and C be an $n \times n$ matrix. Then if B and C are nonsingular matrices, it follows that

$$\text{rank}(BAC) = \text{rank}(BA) = \text{rank}(AC) = \text{rank}(A).$$

By using elementary transformation matrices, any matrix A can be transformed into another matrix of simpler form having the same rank as A .

Theorem 1.11 If A is an $m \times n$ matrix of rank $r > 0$, then nonsingular $m \times m$ and $n \times n$ matrices B and C exist, such that $H = BAC$ and $A = B^{-1}HC^{-1}$, where H is given by

$$\begin{aligned} \text{(a)} \quad & I_r \quad \text{if } r = m = n, & \text{(b)} \quad & \begin{bmatrix} I_r & (0) \end{bmatrix} \quad \text{if } r = m < n, \\ \text{(c)} \quad & \begin{bmatrix} I_r \\ (0) \end{bmatrix} \quad \text{if } r = n < m, & \text{(d)} \quad & \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix} \quad \text{if } r < m, r < n. \end{aligned}$$

Corollary 1.11.1 is an immediate consequence of Theorem 1.11.

Corollary 1.11.1 Let A be an $m \times n$ matrix with $\text{rank}(A) = r > 0$. Then an $m \times r$ matrix F and an $r \times n$ matrix G exist, such that $\text{rank}(F) = \text{rank}(G) = r$ and $A = FG$.

1.10 ORTHOGONAL MATRICES

An $m \times 1$ vector \mathbf{p} is said to be a normalized vector or a unit vector if $\mathbf{p}'\mathbf{p} = 1$. The $m \times 1$ vectors, $\mathbf{p}_1, \dots, \mathbf{p}_n$, where $n \leq m$, are said to be orthogonal if $\mathbf{p}_i'\mathbf{p}_j = 0$ for all $i \neq j$. If in addition, each \mathbf{p}_i is a normalized vector, then the vectors are said to be orthonormal. An $m \times m$ matrix P whose columns form an orthonormal set of vectors is called an orthogonal matrix. It immediately follows that

$$P'P = I_m.$$

Taking the determinant of both sides, we see that

$$|P'P| = |P'| |P| = |P|^2 = |I_m| = 1.$$

Thus, $|P| = +1$ or -1 , so that P is nonsingular, $P^{-1} = P'$, and $PP' = I_m$ in addition to $P'P = I_m$; that is, the rows of P also form an orthonormal set of $m \times 1$ vectors. Some basic properties of orthogonal matrices are summarized in Theorem 1.12.

Theorem 1.12 Let P and Q be $m \times m$ orthogonal matrices and A be any $m \times m$ matrix. Then

- (a) $|P| = \pm 1$,
- (b) $|P'AP| = |A|$,
- (c) PQ is an orthogonal matrix.

One example of an $m \times m$ orthogonal matrix, known as the Helmert matrix, has the form

$$H = \begin{bmatrix} \frac{1}{\sqrt{m}} & \frac{1}{\sqrt{m}} & \frac{1}{\sqrt{m}} & \cdots & \frac{1}{\sqrt{m}} \\ \frac{1}{\sqrt{2}} & -\frac{1}{\sqrt{2}} & 0 & \cdots & 0 \\ \frac{1}{\sqrt{6}} & \frac{1}{\sqrt{6}} & -\frac{2}{\sqrt{6}} & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \frac{1}{\sqrt{m(m-1)}} & \frac{1}{\sqrt{m(m-1)}} & \frac{1}{\sqrt{m(m-1)}} & \cdots & -\frac{(m-1)}{\sqrt{m(m-1)}} \end{bmatrix}.$$

For instance, if $m = 4$, the Helmert matrix is

$$H = \begin{bmatrix} 1/2 & 1/2 & 1/2 & 1/2 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 & 0 \\ 1/\sqrt{6} & 1/\sqrt{6} & -2/\sqrt{6} & 0 \\ 1/\sqrt{12} & 1/\sqrt{12} & 1/\sqrt{12} & -3/\sqrt{12} \end{bmatrix}.$$

Note that if $m \neq n$, it is possible for an $m \times n$ matrix P to satisfy one of the identities, $P'P = I_n$ or $PP' = I_m$, but not both. Such a matrix is sometimes referred to as a semiorthogonal matrix.

An $m \times m$ matrix P is called a permutation matrix if each row and each column of P has a single element 1, while all remaining elements are zeros. As a result, the columns of P will be e_1, \dots, e_m , the columns of I_m , in some order. Note then that the (h, h) th element of $P'P$ will be $e'_i e_i = 1$ for some i , and the (h, l) th element of $P'P$ will be $e'_i e_j = 0$ for some $i \neq j$ if $h \neq l$; that is, a permutation matrix is a special orthogonal matrix. Since there are $m!$ ways of permuting the columns of I_m , there are $m!$ different permutation matrices of order m . If A is also $m \times m$, then PA creates an $m \times m$ matrix by permuting the rows of A , and AP produces a matrix by permuting the columns of A .

1.11 QUADRATIC FORMS

Let \mathbf{x} be an $m \times 1$ vector, \mathbf{y} an $n \times 1$ vector, and A an $m \times n$ matrix. Then the function of \mathbf{x} and \mathbf{y} given by

$$\mathbf{x}'A\mathbf{y} = \sum_{i=1}^m \sum_{j=1}^n x_i y_j a_{ij}$$

is sometimes called a bilinear form in \mathbf{x} and \mathbf{y} . We will be most interested in the special case in which $m = n$, so that A is $m \times m$, and $\mathbf{x} = \mathbf{y}$. In this case, the function above reduces to the function of \mathbf{x} ,

$$f(\mathbf{x}) = \mathbf{x}'A\mathbf{x} = \sum_{i=1}^m \sum_{j=1}^m x_i x_j a_{ij},$$

which is called a quadratic form in \mathbf{x} ; A is referred to as the matrix of the quadratic form. We will always assume that A is a symmetric matrix because, if it is not, A may be replaced by $B = \frac{1}{2}(A + A')$, which is symmetric, without altering $f(\mathbf{x})$; that is,

$$\begin{aligned}\mathbf{x}'B\mathbf{x} &= \frac{1}{2}\mathbf{x}'(A + A')\mathbf{x} = \frac{1}{2}(\mathbf{x}'A\mathbf{x} + \mathbf{x}'A'\mathbf{x}) \\ &= \frac{1}{2}(\mathbf{x}'A\mathbf{x} + \mathbf{x}'A\mathbf{x}) = \mathbf{x}'A\mathbf{x}\end{aligned}$$

because $\mathbf{x}'A'\mathbf{x} = (\mathbf{x}'A'\mathbf{x})' = \mathbf{x}'A\mathbf{x}$. As an example, consider the function

$$f(\mathbf{x}) = x_1^2 + 3x_2^2 + 2x_3^2 + 2x_1x_2 - 2x_2x_3,$$

where \mathbf{x} is 3×1 . The symmetric matrix A satisfying $f(\mathbf{x}) = \mathbf{x}'A\mathbf{x}$ is given by

$$A = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 3 & -1 \\ 0 & -1 & 2 \end{bmatrix}.$$

Every symmetric matrix A and its associated quadratic form is classified into one of the following five categories:

- (a) If $\mathbf{x}'A\mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$, then A is positive definite.
- (b) If $\mathbf{x}'A\mathbf{x} \geq 0$ for all \mathbf{x} and $\mathbf{x}'A\mathbf{x} = 0$ for some $\mathbf{x} \neq \mathbf{0}$, then A is positive semidefinite.
- (c) If $\mathbf{x}'A\mathbf{x} < 0$ for all $\mathbf{x} \neq \mathbf{0}$, then A is negative definite.
- (d) If $\mathbf{x}'A\mathbf{x} \leq 0$ for all \mathbf{x} and $\mathbf{x}'A\mathbf{x} = 0$ for some $\mathbf{x} \neq \mathbf{0}$, then A is negative semidefinite.
- (e) If $\mathbf{x}'A\mathbf{x} > 0$ for some \mathbf{x} and $\mathbf{x}'A\mathbf{x} < 0$ for some \mathbf{x} , then A is indefinite.

Note that the null matrix is actually both positive semidefinite and negative semidefinite.

Positive definite and negative definite matrices are nonsingular, whereas positive semidefinite and negative semidefinite matrices are singular. Sometimes the term *nonnegative definite* will be used to refer to a symmetric matrix that is either positive definite or positive semidefinite. An $m \times m$ matrix B is called a square root of the nonnegative definite $m \times m$ matrix A if $A = BB'$. Sometimes we will denote such a matrix B as $A^{1/2}$. If B is also symmetric, so that $A = B^2$, then B is called the symmetric square root of A .

Quadratic forms play a prominent role in inferential statistics. In Chapter 11, we will develop some of the most important results involving quadratic forms that are of particular interest in statistics.

1.12 COMPLEX MATRICES

Throughout most of this text, we will be dealing with the analysis of vectors and matrices composed of real numbers or variables. However, there are occasions in which an analysis of a real matrix, such as the decomposition of a matrix in the form of a product of other matrices, leads to matrices that contain complex numbers. For this reason, we will briefly summarize in this section some of the basic notation and terminology regarding complex numbers.

Any complex number c can be written in the form

$$c = a + ib,$$

where a and b are real numbers and i represents the imaginary number $\sqrt{-1}$. The real number a is called the real part of c , whereas b is referred to as the imaginary part of c . Thus, the number c is a real number only if b is 0. If we have two complex numbers, $c_1 = a_1 + ib_1$ and $c_2 = a_2 + ib_2$, then their sum is given by

$$c_1 + c_2 = (a_1 + a_2) + i(b_1 + b_2),$$

whereas their product is given by

$$c_1 c_2 = a_1 a_2 - b_1 b_2 + i(a_1 b_2 + a_2 b_1).$$

Corresponding to each complex number $c = a + ib$ is another complex number denoted by \bar{c} and called the complex conjugate of c . The complex conjugate of c is given by $\bar{c} = a - ib$ and satisfies $c\bar{c} = a^2 + b^2$, so that the product of a complex number and its conjugate results in a real number.

A complex number can be represented geometrically by a point in the complex plane, where one of the axes is the real axis and the other axis is the complex or imaginary axis. Thus, the complex number $c = a + ib$ would be represented by the point (a, b) in this complex plane. Alternatively, we can use the polar coordinates (r, θ) , where r is the length of the line from the origin to the point (a, b) and θ is the angle between this line and the positive half of the real axis. The relationship between a and b , and r and θ is then given by

$$a = r \cos(\theta), \quad b = r \sin(\theta).$$

Writing c in terms of the polar coordinates, we have

$$c = r \cos(\theta) + ir \sin(\theta),$$

or, after using Euler's formula, simply $c = re^{i\theta}$. The absolute value, also sometimes called the modulus, of the complex number c is defined to be r . This is, of course, always a nonnegative real number, and because $a^2 + b^2 = r^2$, we have

$$|c| = |a + ib| = \sqrt{a^2 + b^2}.$$

We also find that

$$\begin{aligned} |c_1 c_2| &= \sqrt{(a_1 a_2 - b_1 b_2)^2 + (a_1 b_2 + a_2 b_1)^2} \\ &= \sqrt{(a_1^2 + b_1^2)(a_2^2 + b_2^2)} = |c_1| |c_2|. \end{aligned}$$

Using this identity repeatedly, we also see that for any complex number c and any positive integer n , $|c^n| = |c|^n$.

A useful identity relating a complex number c and its conjugate to the absolute value of c is

$$c\bar{c} = |c|^2.$$

Applying this identity to the sum of two complex numbers $c_1 + c_2$ and noting that $c_1\bar{c}_2 + \bar{c}_1 c_2 \leq 2|c_1||c_2|$, we get

$$\begin{aligned} |c_1 + c_2|^2 &= (c_1 + c_2)\overline{(c_1 + c_2)} = (c_1 + c_2)(\bar{c}_1 + \bar{c}_2) \\ &= c_1\bar{c}_1 + c_1\bar{c}_2 + c_2\bar{c}_1 + c_2\bar{c}_2 \\ &\leq |c_1|^2 + 2|c_1||c_2| + |c_2|^2 \\ &= (|c_1| + |c_2|)^2. \end{aligned}$$

From this result, we get the important inequality, $|c_1 + c_2| \leq |c_1| + |c_2|$, known as the triangle inequality.

A complex matrix is simply a matrix whose elements are complex numbers. As a result, a complex matrix can be written as the sum of a real matrix and an imaginary matrix; that is, if C is an $m \times n$ complex matrix then it can be expressed as

$$C = A + iB,$$

where both A and B are $m \times n$ real matrices. The complex conjugate of C , denoted \bar{C} , is simply the matrix containing the complex conjugates of the elements of C ; that is,

$$\bar{C} = A - iB.$$

The conjugate transpose of C is $C^* = \bar{C}'$. If the complex matrix C is square and $C^* = C$, so that $c_{ij} = \bar{c}_{ji}$, then C is said to be Hermitian. Note that if C is Hermitian and C is a real matrix, then C is symmetric. The $m \times m$ matrix C is said to be unitary if $C^*C = I_m$, which is the generalization of the concept of orthogonal matrices to complex matrices because if C is real, then $C^* = C'$.

1.13 RANDOM VECTORS AND SOME RELATED STATISTICAL CONCEPTS

In this section, we review some of the basic definitions and results in distribution theory that will be needed later in this text. A more comprehensive treatment of this

subject can be found in books on statistical theory such as Casella and Berger (2002) or Lindgren (1993). To be consistent with our notation in which we use a capital letter to denote a matrix, a bold lowercase letter for a vector, and a lowercase letter for a scalar, we will use a lowercase letter instead of the more conventional capital letter to denote a scalar random variable.

A random variable x is said to be discrete if its collection of possible values, R_x , is a countable set. In this case, x has a probability function $p_x(t)$ satisfying $p_x(t) = P(x = t)$, for $t \in R_x$, and $p_x(t) = 0$, for $t \notin R_x$. A continuous random variable x , on the other hand, has for its range, R_x , an uncountably infinite set. Associated with each continuous random variable x is a density function $f_x(t)$ satisfying $f_x(t) > 0$, for $t \in R_x$, and $f_x(t) = 0$, for $t \notin R_x$. Probabilities for x are obtained by integration; if B is a subset of the real line, then

$$P(x \in B) = \int_B f_x(t) dt.$$

For both discrete and continuous x , we have $P(x \in R_x) = 1$.

The expected value of a real-valued function of x , $g(x)$, gives the average observed value of $g(x)$. This expectation, denoted $E[g(x)]$, is given by

$$E[g(x)] = \sum_{t \in R_x} g(t)p_x(t),$$

if x is discrete, and

$$E[g(x)] = \int_{-\infty}^{\infty} g(t)f_x(t) dt,$$

if x is continuous. Properties of the expectation operator follow directly from properties of sums and integrals. For instance, if x is a random variable and α and β are constants, then the expectation operator satisfies the properties

$$E(\alpha) = \alpha$$

and

$$E[\alpha g_1(x) + \beta g_2(x)] = \alpha E[g_1(x)] + \beta E[g_2(x)],$$

where g_1 and g_2 are any real-valued functions. The expected values of a random variable x given by $E(x^k)$, $k = 1, 2, \dots$ are known as the moments of x . These moments are important for both descriptive and theoretical purposes. The first few moments can be used to describe certain features of the distribution of x . For instance, the first moment or mean of x , $\mu_x = E(x)$, locates a central value of the distribution. The variance of x , denoted σ_x^2 or $\text{var}(x)$, is defined as

$$\sigma_x^2 = \text{var}(x) = E[(x - \mu_x)^2] = E(x^2) - \mu_x^2,$$

so that it is a function of the first and second moments of x . The variance gives a measure of the dispersion of the observed values of x about the central value μ_x . Using properties of expectation, it is easily verified that

$$\text{var}(\alpha + \beta x) = \beta^2 \text{var}(x).$$

All of the moments of a random variable x are embedded in a function called the moment generating function of x . This function is defined as a particular expectation; specifically, the moment generating function of x , $m_x(t)$, is given by

$$m_x(t) = E(e^{tx}),$$

provided this expectation exists for values of t in a neighborhood of 0. Otherwise, the moment generating function does not exist. If the moment generating function of x does exist, then we can obtain any moment from it because

$$\left. \frac{d^k}{dt^k} m_x(t) \right|_{t=0} = E(x^k).$$

More importantly, the moment generating function characterizes the distribution of x in that, under certain conditions, no two different distributions have the same moment generating function.

We now focus on some particular families of distributions that we will encounter later in this text. A random variable x is said to have a univariate normal distribution with mean μ and variance σ^2 , indicated by $x \sim N(\mu, \sigma^2)$, if the density of x is given by

$$f_x(t) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(t-\mu)^2/2\sigma^2}, \quad -\infty < t < \infty.$$

The corresponding moment generating function is

$$m_x(t) = e^{\mu t + \sigma^2 t^2/2}.$$

A special member of this family of normal distributions is the standard normal distribution $N(0, 1)$. The importance of this distribution follows from the fact that if $x \sim N(\mu, \sigma^2)$, then the standardizing transformation $z = (x - \mu)/\sigma$ yields a random variable z that has the standard normal distribution. By differentiating the moment generating function of $z \sim N(0, 1)$, it is easy to verify that the first six moments of z , which we will need in Chapter 11, are 0, 1, 0, 3, 0, and 15, respectively.

If r is a positive integer, then a random variable v has a chi-squared distribution with r degrees of freedom, written $v \sim \chi_r^2$, if its density function is

$$f_v(t) = \frac{t^{(r/2)-1} e^{-t/2}}{2^{r/2} \Gamma(r/2)}, \quad t > 0,$$

where $\Gamma(r/2)$ is the gamma function evaluated at $r/2$. The moment generating function of v is given by $m_v(t) = (1 - 2t)^{-r/2}$, for $t < \frac{1}{2}$. The importance of

the chi-squared distribution arises from its connection to the normal distribution. If $z \sim N(0, 1)$, then $z^2 \sim \chi_1^2$. Further, if z_1, \dots, z_r are independent random variables with $z_i \sim N(0, 1)$ for $i = 1, \dots, r$, then

$$\sum_{i=1}^r z_i^2 \sim \chi_r^2. \quad (1.5)$$

The chi-squared distribution mentioned above is sometimes referred to as a central chi-squared distribution because it is actually a special case of a more general family of distributions known as the noncentral chi-squared distributions. These noncentral chi-squared distributions are also related to the normal distribution. If x_1, \dots, x_r are independent random variables with $x_i \sim N(\mu_i, 1)$, then

$$\sum_{i=1}^r x_i^2 \sim \chi_r^2(\lambda), \quad (1.6)$$

where $\chi_r^2(\lambda)$ denotes the noncentral chi-squared distribution with r degrees of freedom and noncentrality parameter

$$\lambda = \frac{1}{2} \sum_{i=1}^r \mu_i^2;$$

that is, the noncentral chi-squared density, which we will not give here, depends not only on the parameter r but also on the parameter λ . Since (1.6) reduces to (1.5) when $\mu_i = 0$ for all i , we see that the distribution $\chi_r^2(\lambda)$ corresponds to the central chi-squared distribution χ_r^2 when $\lambda = 0$.

A distribution related to the chi-squared distribution is the F distribution with r_1 and r_2 degrees of freedom, denoted by F_{r_1, r_2} . If $y \sim F_{r_1, r_2}$, then the density function of y is

$$f_y(t) = \frac{\Gamma\{(r_1 + r_2)/2\}}{\Gamma(r_1/2)\Gamma(r_2/2)} \left(\frac{r_1}{r_2}\right)^{r_1/2} t^{(r_1-2)/2} \left(1 + \frac{r_1}{r_2}t\right)^{-(r_1+r_2)/2}, \quad t > 0.$$

The importance of this distribution arises from the fact that if v_1 and v_2 are independent random variables with $v_1 \sim \chi_{r_1}^2$ and $v_2 \sim \chi_{r_2}^2$, then the ratio

$$t = \frac{v_1/r_1}{v_2/r_2}$$

has the F distribution with r_1 and r_2 degrees of freedom.

The concept of a random variable can be extended to that of a random vector. A sequence of related random variables x_1, \dots, x_m is modeled by a joint or multivariate probability function $p_{\mathbf{x}}(\mathbf{t})$ if all of the random variables are discrete, and a multivariate density function $f_{\mathbf{x}}(\mathbf{t})$ if all of the random variables are continuous, where $\mathbf{x} = (x_1, \dots, x_m)'$ and $\mathbf{t} = (t_1, \dots, t_m)'$. For instance, if they are continuous and B is a region in R^m , then the probability that \mathbf{x} falls in B is

$$P(\mathbf{x} \in B) = \int \cdots \int_B f_{\mathbf{x}}(\mathbf{t}) dt_1 \cdots dt_m,$$

whereas the expected value of the real-valued function $g(\mathbf{x})$ of \mathbf{x} is given by

$$E[g(\mathbf{x})] = \int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} g(\mathbf{t}) f_{\mathbf{x}}(\mathbf{t}) dt_1 \cdots dt_m.$$

The random variables x_1, \dots, x_m are said to be independent, a concept we have already referred to, if and only if the joint probability function or density function factors into the product of the marginal probability or density functions; that is, in the continuous case, x_1, \dots, x_m are independent if and only if

$$f_{\mathbf{x}}(\mathbf{t}) = f_{x_1}(t_1) \cdots f_{x_m}(t_m),$$

for all \mathbf{t} .

The mean vector of \mathbf{x} , denoted $\boldsymbol{\mu}$, is the vector of expected values of the x_i 's; that is,

$$\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)' = E(\mathbf{x}) = [E(x_1), \dots, E(x_m)]'.$$

A measure of the linear relationship between x_i and x_j is given by the covariance of x_i and x_j , which is denoted $\text{cov}(x_i, x_j)$ or σ_{ij} and is defined by

$$\sigma_{ij} = \text{cov}(x_i, x_j) = E[(x_i - \mu_i)(x_j - \mu_j)] = E(x_i x_j) - \mu_i \mu_j. \quad (1.7)$$

When $i = j$, this covariance reduces to the variance of x_i ; that is, $\sigma_{ii} = \sigma_i^2 = \text{var}(x_i)$. When $i \neq j$ and x_i and x_j are independent, then $\text{cov}(x_i, x_j) = 0$ because in this case $E(x_i x_j) = \mu_i \mu_j$. If $\alpha_1, \alpha_2, \beta_1$ and β_2 are constants, then

$$\text{cov}(\alpha_1 + \beta_1 x_i, \alpha_2 + \beta_2 x_j) = \beta_1 \beta_2 \text{cov}(x_i, x_j).$$

The matrix Ω , which has σ_{ij} as its (i, j) th element, is called the variance–covariance matrix, or simply the covariance matrix, of \mathbf{x} . This matrix will be also denoted sometimes by $\text{var}(\mathbf{x})$ or $\text{cov}(\mathbf{x}, \mathbf{x})$. Clearly, $\sigma_{ij} = \sigma_{ji}$ so that Ω is a symmetric matrix. Using (1.7), we obtain the matrix formulation for Ω ,

$$\Omega = \text{var}(\mathbf{x}) = E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})'] = E(\mathbf{x}\mathbf{x}') - \boldsymbol{\mu}\boldsymbol{\mu}'.$$

If $\boldsymbol{\alpha}$ is an $m \times 1$ vector of constants and we define the random variable $y = \boldsymbol{\alpha}'\mathbf{x}$, then

$$\begin{aligned} E(y) &= E(\boldsymbol{\alpha}'\mathbf{x}) = E\left(\sum_{i=1}^m \alpha_i x_i\right) = \sum_{i=1}^m \alpha_i E(x_i) \\ &= \sum_{i=1}^m \alpha_i \mu_i = \boldsymbol{\alpha}'\boldsymbol{\mu}. \end{aligned}$$

If, in addition, β is another $m \times 1$ vector of constants and $w = \beta'x$, then

$$\begin{aligned} \text{cov}(y, w) &= \text{cov}(\alpha'x, \beta'x) = \text{cov}\left(\sum_{i=1}^m \alpha_i x_i, \sum_{j=1}^m \beta_j x_j\right) \\ &= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \beta_j \text{cov}(x_i, x_j) = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \beta_j \sigma_{ij} = \alpha' \Omega \beta. \end{aligned}$$

In particular, $\text{var}(y) = \text{cov}(y, y) = \alpha' \Omega \alpha$. Because this holds for any choice of α and because the variance is always nonnegative, Ω must be a nonnegative definite matrix. More generally, if A is a $p \times m$ matrix of constants and $y = Ax$, then

$$E(y) = E(Ax) = AE(x) = A\mu, \quad (1.8)$$

$$\begin{aligned} \text{var}(y) &= E[\{y - E(y)\}\{y - E(y)\}'] = E[(Ax - A\mu)(Ax - A\mu)'] \\ &= E[A(x - \mu)(x - \mu)'A'] = A\{E[(x - \mu)(x - \mu)']\}A' \\ &= A\Omega A'. \end{aligned} \quad (1.9)$$

Thus, the mean vector and covariance matrix of the transformed vector, Ax , is $A\mu$ and $A\Omega A'$. If v and w are random vectors, then the matrix of covariances between components of v and components of w is given by

$$\text{cov}(v, w) = E(vw') - E(v)E(w)'.$$

In particular, if $v = Ax$ and $w = Bx$, then

$$\text{cov}(v, w) = A \text{cov}(x, x) B' = A \text{var}(x) B' = A\Omega B'.$$

A measure of the linear relationship between x_i and x_j that is unaffected by the measurement scales of x_i and x_j is called the correlation. We denote this by ρ_{ij} and it is defined as

$$\rho_{ij} = \frac{\text{cov}(x_i, x_j)}{\sqrt{\text{var}(x_i) \text{var}(x_j)}} = \frac{\sigma_{ij}}{\sqrt{\sigma_{ii} \sigma_{jj}}}.$$

When $i = j$, $\rho_{ij} = 1$. The correlation matrix P , which has ρ_{ij} as its (i, j) th element, can be expressed in terms of the corresponding covariance matrix Ω and the diagonal matrix $D_\Omega^{-1/2} = \text{diag}(\sigma_{11}^{-1/2}, \dots, \sigma_{mm}^{-1/2})$; specifically,

$$P = D_\Omega^{-1/2} \Omega D_\Omega^{-1/2}. \quad (1.10)$$

For any $m \times 1$ vector α , we have

$$\alpha' P \alpha = \alpha' D_\Omega^{-1/2} \Omega D_\Omega^{-1/2} \alpha = \beta' \Omega \beta,$$

where $\beta = D_\Omega^{-1/2}\alpha$, and so P must be nonnegative definite because Ω is. In particular, if e_i is the i th column of the $m \times m$ identity matrix, then

$$\begin{aligned} (e_i + e_j)'P(e_i + e_j) &= (P)_{ii} + (P)_{ij} + (P)_{ji} + (P)_{jj} \\ &= 2(1 + \rho_{ij}) \geq 0 \end{aligned}$$

and

$$\begin{aligned} (e_i - e_j)'P(e_i - e_j) &= (P)_{ii} - (P)_{ij} - (P)_{ji} + (P)_{jj} \\ &= 2(1 - \rho_{ij}) \geq 0, \end{aligned}$$

from which we obtain the inequality, $-1 \leq \rho_{ij} \leq 1$.

Typically, means, variances, and covariances are unknown and so they must be estimated from a sample. Suppose x_1, \dots, x_n represents a random sample of a random variable x that has some distribution with mean μ and variance σ^2 . These quantities can be estimated by the sample mean and the sample variance given by

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i, \\ s^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n x_i^2 - n\bar{x}^2 \right). \end{aligned}$$

In the multivariate setting, we have analogous estimators for μ and Ω ; if x_1, \dots, x_n is a random sample of an $m \times 1$ random vector x having mean vector μ and covariance matrix Ω , then the sample mean vector and sample covariance matrix are given by

$$\begin{aligned} \bar{x} &= \frac{1}{n} \sum_{i=1}^n x_i, \\ S &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})' = \frac{1}{n-1} \left(\sum_{i=1}^n x_i x_i' - n\bar{x}\bar{x}' \right). \end{aligned}$$

The sample covariance matrix can be then used in (1.10) to obtain an estimator of the correlation matrix, P ; that is, if we define the diagonal matrix $D_S^{-1/2} = \text{diag}(s_{11}^{-1/2}, \dots, s_{mm}^{-1/2})$, then the correlation matrix can be estimated by the sample correlation matrix defined as

$$R = D_S^{-1/2} S D_S^{-1/2}.$$

One particular joint distribution that we will consider is the multivariate normal distribution. This distribution can be defined in terms of independent standard normal

random variables. Let z_1, \dots, z_m be independently distributed as $N(0, 1)$, and put $\mathbf{z} = (z_1, \dots, z_m)'$. The density function of \mathbf{z} is then given by

$$f(\mathbf{z}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z_i^2\right) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\mathbf{z}'\mathbf{z}\right).$$

Because $E(\mathbf{z}) = \mathbf{0}$ and $\text{var}(\mathbf{z}) = I_m$, this particular m -dimensional multivariate normal distribution, known as the standard multivariate normal distribution, is denoted as $N_m(\mathbf{0}, I_m)$. If $\boldsymbol{\mu}$ is an $m \times 1$ vector of constants and T is an $m \times m$ nonsingular matrix, then $\mathbf{x} = \boldsymbol{\mu} + T\mathbf{z}$ is said to have the m -dimensional multivariate normal distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\Omega = TT'$. This is indicated by $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$. For instance, if $m = 2$, the vector $\mathbf{x} = (x_1, x_2)'$ has a bivariate normal distribution and its density, induced by the transformation $\mathbf{x} = \boldsymbol{\mu} + T\mathbf{z}$, can be shown to be

$$f(\mathbf{x}) = \frac{1}{2\pi\sqrt{\sigma_{11}\sigma_{22}(1-\rho^2)}} \exp\left(-\frac{1}{2(1-\rho^2)}\left\{\frac{(x_1-\mu_1)^2}{\sigma_{11}} - 2\rho\left(\frac{x_1-\mu_1}{\sqrt{\sigma_{11}}}\right)\left(\frac{x_2-\mu_2}{\sqrt{\sigma_{22}}}\right) + \frac{(x_2-\mu_2)^2}{\sigma_{22}}\right\}\right), \quad (1.11)$$

for all $\mathbf{x} \in R^2$, where $\rho = \rho_{12}$ is the correlation coefficient. When $\rho = 0$, this density factors into the product of the marginal densities, so x_1 and x_2 are independent if and only if $\rho = 0$. The cumbersome-looking density function given in (1.11) can be more conveniently expressed by using matrix notation. It is straightforward to verify that this density is identical to

$$f(\mathbf{x}) = \frac{1}{2\pi|\Omega|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}. \quad (1.12)$$

The density function of an m -variate normal random vector is very similar to the function given in (1.12). If $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, then its density is

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{m/2}|\Omega|^{1/2}} \exp\left\{-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\right\}, \quad (1.13)$$

for all $\mathbf{x} \in R^m$.

If Ω is positive semidefinite, then $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ is said to have a singular normal distribution. In this case, Ω^{-1} does not exist and so the multivariate normal density cannot be written in the form given in (1.13). However, the random vector \mathbf{x} can still be expressed in terms of independent standard normal random variables. Suppose that $\text{rank}(\Omega) = r$ and U is an $m \times r$ matrix satisfying $UU' = \Omega$. Then $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ if \mathbf{x} is distributed the same as $\boldsymbol{\mu} + U\mathbf{z}$, where now $\mathbf{z} \sim N_r(\mathbf{0}, I_r)$.

An important property of the multivariate normal distribution is that a linear transformation of a multivariate normal vector yields a multivariate normal vector; that is, if $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ and A is a $p \times m$ matrix of constants, then $\mathbf{y} = A\mathbf{x}$ has a p -variate normal distribution. In particular, from (1.8) and (1.9), we know that $\mathbf{y} \sim N_p(A\boldsymbol{\mu}, A\Omega A')$.

We next consider spherical and elliptical distributions that are extensions of multivariate normal distributions. In particular, a spherical distribution is an extension of the standard multivariate normal distribution $N_m(\mathbf{0}, I_m)$, whereas an elliptical distribution is an extension of the multivariate normal distribution $N_m(\boldsymbol{\mu}, \Omega)$. An $m \times 1$ random vector \mathbf{x} has a spherical distribution if \mathbf{x} and $P\mathbf{x}$ have the same distribution for all $m \times m$ orthogonal matrices P . If \mathbf{x} has a spherical distribution with a density function, then this density function depends on \mathbf{x} only through the value of $\mathbf{x}'\mathbf{x}$; that is, the density function of \mathbf{x} can be written as $g(\mathbf{x}'\mathbf{x})$ for some function g . The term *spherical* distribution then arises from the fact that the density function is the same for all points \mathbf{x} that lie on the sphere $\mathbf{x}'\mathbf{x} = c$, where c is a nonnegative constant. Clearly $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$ has a spherical distribution because for any $m \times m$ orthogonal matrix P , $P\mathbf{z} \sim N_m(\mathbf{0}, I_m)$. An example of a nonnormal spherical distribution is the uniform distribution; that is, if \mathbf{u} is a randomly selected point on the surface of the unit sphere in R^m , then \mathbf{u} has a spherical distribution. In fact, if the $m \times 1$ random vector \mathbf{x} has a spherical distribution, then it can be expressed as

$$\mathbf{x} = w\mathbf{u}, \quad (1.14)$$

where \mathbf{u} is uniformly distributed on the m -dimensional unit sphere, w is a nonnegative random variable, and \mathbf{u} and w are independently distributed. It is easy to verify that when \mathbf{z} has the distribution $N_m(\mathbf{0}, I_m)$, then (1.14) takes the form

$$\mathbf{z} = v\mathbf{u},$$

where $v^2 \sim \chi_m^2$. Thus, if the $m \times 1$ random vector \mathbf{x} has a spherical distribution, then it can also be expressed as

$$\mathbf{x} = w\mathbf{u} = wv^{-1}\mathbf{z} = s\mathbf{z},$$

where again \mathbf{z} has the distribution $N_m(\mathbf{0}, I_m)$, $s = wv^{-1}$ is a nonnegative random variable, and \mathbf{z} and s are independently distributed. The contaminated normal distributions and the multivariate t distributions are other examples of spherical distributions. A random vector \mathbf{x} having a contaminated normal distribution can be expressed as $\mathbf{x} = s\mathbf{z}$, where $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$ independently of s , which takes on the values σ and 1 with probabilities p and $1 - p$, respectively, and $\sigma \neq 1$ is a positive constant. If $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$ independently of $v^2 \sim \chi_n^2$, then the random vector $\mathbf{x} = n^{1/2}\mathbf{z}/v$ has a multivariate t distribution with n degrees of freedom.

We generalize from spherical distributions to elliptical distributions in the same way that $N_m(\mathbf{0}, I_m)$ was generalized to $N_m(\boldsymbol{\mu}, \Omega)$. An $m \times 1$ random vector \mathbf{y} has an elliptical distribution with parameters $\boldsymbol{\mu}$ and Ω if it can be expressed as

$$\mathbf{y} = \boldsymbol{\mu} + T\mathbf{x},$$

where T is $m \times r$, $TT' = \Omega$, $\text{rank}(\Omega) = r$, and the $r \times 1$ random vector \mathbf{x} has a spherical distribution. Using (1.14), we then have

$$\mathbf{y} = \boldsymbol{\mu} + wT\mathbf{u},$$

where the random variable $w \geq 0$ is independent of \mathbf{u} , which is uniformly distributed on the r -dimensional unit sphere. If Ω is nonsingular and \mathbf{y} has a density, then it

depends on \mathbf{y} only through the value of $(\mathbf{y} - \boldsymbol{\mu})' \Omega^{-1} (\mathbf{y} - \boldsymbol{\mu})$; that is, the density is the same for all points \mathbf{y} that lie on the ellipsoid $(\mathbf{y} - \boldsymbol{\mu})' \Omega^{-1} (\mathbf{y} - \boldsymbol{\mu}) = c$, where c is a nonnegative constant. A more detailed discussion about spherical and elliptical distributions can be found in Fang et al. (1990).

One of the most widely used procedures in statistics is regression analysis. We will briefly describe this analysis here and later use regression analysis to illustrate some of the matrix methods developed in this text. Some good references on regression are Kutner et al. (2005), Rencher and Schaalje (2008), and Sen and Srivastava (1990). In the typical regression problem, one wishes to study the relationship between some response variable, say y , and k explanatory variables x_1, \dots, x_k . For instance, y might be the yield of some product of a manufacturing process, whereas the explanatory variables are conditions affecting the production process, such as temperature, humidity, pressure, and so on. A model relating the x_j 's to y is given by

$$y = \beta_0 + \beta_1 x_1 + \dots + \beta_k x_k + \epsilon, \quad (1.15)$$

where β_0, \dots, β_k are unknown parameters and ϵ is a random error, that is, a random variable, with $E(\epsilon) = 0$. In what is known as ordinary least squares regression, we also have the errors as independent random variables with common variance σ^2 ; that is, if ϵ_i and ϵ_j are random errors associated with the responses y_i and y_j , then $\text{var}(\epsilon_i) = \text{var}(\epsilon_j) = \sigma^2$ and $\text{cov}(\epsilon_i, \epsilon_j) = 0$. The model given in (1.15) is an example of a linear model because it is a linear function of the parameters. It need not be linear in the x_j 's so that, for instance, we might have $x_2 = x_1^2$. Because the parameters are unknown, they must be estimated and this will be possible if we have some observed values of y and the corresponding x_j 's. Thus, for the i th observation, suppose that the explanatory variables are set to the values x_{i1}, \dots, x_{ik} yielding the response y_i , and this is done for $i = 1, \dots, N$, where $N > k + 1$. If model (1.15) holds, then we should have, approximately,

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_k x_{ik}$$

for each i . This can be written as the matrix equation

$$\mathbf{y} = X\boldsymbol{\beta}$$

if we define

$$\mathbf{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{bmatrix}, \quad X = \begin{bmatrix} 1 & x_{11} & \dots & x_{1k} \\ 1 & x_{21} & \dots & x_{2k} \\ \vdots & \vdots & \dots & \vdots \\ 1 & x_{N1} & \dots & x_{Nk} \end{bmatrix}.$$

One method of estimating the β_j 's, which we will discuss from time to time in this text, is called the method of least squares. If $\hat{\boldsymbol{\beta}} = (\hat{\beta}_0, \dots, \hat{\beta}_k)'$ is an estimate of the

parameter vector β , then $\hat{\mathbf{y}} = X\hat{\beta}$ is the vector of fitted values, whereas $\mathbf{y} - \hat{\mathbf{y}}$ gives the vector of errors or deviations of the actual responses from the corresponding fitted values, and

$$f(\hat{\beta}) = (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta})$$

gives the sum of squares of these errors. The method of least squares selects as $\hat{\beta}$ any vector that minimizes the function $f(\hat{\beta})$. We will see later that any such vector satisfies the system of linear equations, sometimes referred to as the normal equations,

$$X'X\hat{\beta} = X'\mathbf{y}.$$

If X has full column rank, that is, $\text{rank}(X) = k + 1$, then $(X'X)^{-1}$ exists and so the least squares estimator of β is unique and is given by

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

PROBLEMS

- 1.1** Show that the scalar properties $ab = 0$ implies $a = 0$ or $b = 0$, and $ab = ac$ for $a \neq 0$ implies that $b = c$ do not extend to matrices by finding
 - (a) 2×2 nonnull matrices A and B for which $AB = (0)$,
 - (b) 2×2 matrices A , B , and C , with A being nonnull, such that $AB = AC$, yet $B \neq C$.
- 1.2** Let A be an $m \times m$ idempotent matrix. Show that
 - (a) $I_m - A$ is idempotent,
 - (b) BAB^{-1} is idempotent, where B is any $m \times m$ nonsingular matrix.
- 1.3** Let A and B be $m \times m$ symmetric matrices. Show that AB is symmetric if and only if $AB = BA$.
- 1.4** Prove Theorem 1.3(e); that is, if A is an $m \times n$ matrix, show that $\text{tr}(A'A) = 0$ if and only if $A = (0)$.
- 1.5** Show that
 - (a) if \mathbf{x} and \mathbf{y} are $m \times 1$ vectors, $\text{tr}(\mathbf{x}\mathbf{y}') = \mathbf{x}'\mathbf{y}$,
 - (b) if A and B are $m \times m$ matrices and B is nonsingular, $\text{tr}(BAB^{-1}) = \text{tr}(A)$.
- 1.6** Suppose A is $m \times n$ and B is $n \times m$. Show that $\text{tr}(AB) = \text{tr}(A'B')$.
- 1.7** Suppose that A , B , and C are $m \times m$ matrices. Show that if they are symmetric matrices, then $\text{tr}(ABC) = \text{tr}(ACB)$.
- 1.8** Prove Theorem 1.4.
- 1.9** Show that any square matrix can be written as the sum of a symmetric matrix and a skew-symmetric matrix.
- 1.10** Let A and B be $m \times m$ symmetric matrices. Show that $AB - BA$ is a skew-symmetric matrix.

1.11 Suppose that A is an $m \times m$ skew-symmetric matrix. Show that $-A^2$ is a non-negative definite matrix.

1.12 Define the $m \times m$ matrices A , B , and C as

$$A = \begin{bmatrix} b_{11} + c_{11} & b_{12} + c_{12} & \cdots & b_{1m} + c_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{bmatrix},$$

$$B = \begin{bmatrix} b_{11} & b_{12} & \cdots & b_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{bmatrix},$$

$$C = \begin{bmatrix} c_{11} & c_{12} & \cdots & c_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{bmatrix}.$$

Prove that $|A| = |B| + |C|$.

1.13 Verify the results of Theorem 1.8.

1.14 Suppose that A and B are $m \times m$ nonnull matrices satisfying $AB = (0)$. Show that both A and B must be singular matrices.

1.15 Consider the 4×4 matrix

$$A = \begin{bmatrix} 1 & 2 & 1 & 1 \\ 0 & 1 & 2 & 0 \\ 1 & 2 & 2 & 1 \\ 0 & -1 & 1 & 2 \end{bmatrix}.$$

Find the determinant of A by using the cofactor expansion formula on the first column of A .

1.16 Using the matrix A from the previous problem, verify (1.3) when $i = 1$ and $k = 2$.

1.17 Prove Theorem 1.6.

1.18 Let λ be a variable, and consider the determinant of $A - \lambda I_m$, where A is an $m \times m$ matrix, as a function of λ . What type of function of λ is this?

1.19 Find the adjoint matrix of the matrix A given in Problem 1.15. Use this to obtain the inverse of A .

1.20 Using elementary transformations, determine matrices B and C so that $BAC = I_4$ for the matrix A given in Problem 1.15. Use B and C to compute the inverse of A ; that is, take the inverse of both sides of the equation $BAC = I_4$ and then solve for A^{-1} .

1.21 Compute the inverse of

(a) $I_m + \mathbf{1}_m \mathbf{1}_m'$,

(b) $I_m + e_1 \mathbf{1}_m'$.

1.22 Show that

(a) the determinant of a triangular matrix is the product of its diagonal elements,

(b) the inverse of a lower triangular matrix is a lower triangular matrix.

1.23 Let \mathbf{a} and \mathbf{b} be $m \times 1$ vectors and D be an $m \times m$ diagonal matrix. Use Corollary 1.9.2 to find an expression for the inverse of $D + \alpha \mathbf{a} \mathbf{b}'$, where α is a scalar.

1.24 Let $A_{\#}$ be the adjoint matrix of an $m \times m$ matrix A . Show that

(a) $|A_{\#}| = |A|^{m-1}$,

(b) $(\alpha A)_{\#} = \alpha^{m-1} A_{\#}$, where α is a scalar.

1.25 Consider the $m \times m$ partitioned matrix

$$A = \begin{bmatrix} A_{11} & (0) \\ A_{21} & A_{22} \end{bmatrix},$$

where the $m_1 \times m_1$ matrix A_{11} and the $m_2 \times m_2$ matrix A_{22} are nonsingular. Obtain an expression for A^{-1} in terms of A_{11} , A_{22} , and A_{21} .

1.26 Let

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{12}' & A_{22} \end{bmatrix},$$

where A_{11} is $m_1 \times m_1$, A_{22} is $m_2 \times m_2$, and A_{12} is $m_1 \times m_2$. Show that if A is positive definite, then A_{11} and A_{22} are also positive definite.

1.27 Find the rank of the 4×4 matrix

$$A = \begin{bmatrix} 2 & 0 & 1 & -1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & 2 & 0 \\ 2 & 0 & 0 & -2 \end{bmatrix}.$$

1.28 Use elementary transformations to transform the matrix A given in Problem 1.27 to a matrix H having the form given in Theorem 1.11. Consequently, determine matrices B and C so that $BAC = H$.

1.29 Prove parts (b) and (c) of Theorem 1.12.

1.30 List all permutation matrices of order 3.

1.31 Consider the 3×3 matrix

$$P = \frac{1}{\sqrt{6}} \begin{bmatrix} \sqrt{2} & \sqrt{2} & \sqrt{2} \\ \sqrt{3} & -\sqrt{3} & 0 \\ p_{31} & p_{32} & p_{33} \end{bmatrix}.$$

Find values for p_{31} , p_{32} , and p_{33} so that P is an orthogonal matrix. Is your solution unique?

- 1.32** Give the conditions on the $m \times 1$ vector \mathbf{x} so that the matrix $H = I_m - 2\mathbf{x}\mathbf{x}'$ is orthogonal.
- 1.33** Suppose the $m \times m$ orthogonal matrix P is partitioned as $P = [P_1 \ P_2]$, where P_1 is $m \times m_1$, P_2 is $m \times m_2$, and $m_1 + m_2 = m$. Show that $P_1'P_1 = I_{m_1}$, $P_2'P_2 = I_{m_2}$, and $P_1P_1' + P_2P_2' = I_m$.
- 1.34** Let A , B , and C be $m \times n$, $n \times p$, and $n \times n$ matrices, respectively, while \mathbf{x} is an $n \times 1$ vector. Show that
- (a) $A\mathbf{x} = \mathbf{0}$ for all choices of \mathbf{x} if and only if $A = (0)$,
 - (b) $A\mathbf{x} = \mathbf{0}$ if and only if $A'A\mathbf{x} = \mathbf{0}$,
 - (c) $A = (0)$ if $A'A = (0)$,
 - (d) $AB = (0)$ if and only if $A'AB = (0)$,
 - (e) $\mathbf{x}'C\mathbf{x} = 0$ for all \mathbf{x} if and only if $C' = -C$.
- 1.35** For each of the following, find the 3×3 symmetric matrix A so that the given identity holds:
- (a) $\mathbf{x}'A\mathbf{x} = x_1^2 + 2x_2^2 - x_3^2 + 4x_1x_2 - 6x_1x_3 + 8x_2x_3$.
 - (b) $\mathbf{x}'A\mathbf{x} = 3x_1^2 + 5x_2^2 + 2x_3^2 + 2x_1x_2 + 2x_1x_3 + 4x_2x_3$.
 - (c) $\mathbf{x}'A\mathbf{x} = 2x_1x_2 + 2x_1x_3 + 2x_2x_3$.
- 1.36** Let \mathbf{x} be a 4×1 vector. Find symmetric matrices A_1 and A_2 such that

$$\begin{aligned}\mathbf{x}'A_1\mathbf{x} &= (x_1 + x_2 - 2x_3)^2 + (x_3 - x_4)^2, \\ \mathbf{x}'A_2\mathbf{x} &= (x_1 - x_2 - x_3)^2 + (x_1 + x_2 - x_4)^2.\end{aligned}$$

- 1.37** Let A be an $m \times m$ matrix, and suppose that a real $n \times m$ matrix T exists such that $T'T = A$. Show that A must be nonnegative definite.
- 1.38** Prove that a nonnegative definite matrix must have nonnegative diagonal elements; that is, show that if a symmetric matrix has any negative diagonal elements, then it is not nonnegative definite. Show that the converse is not true; that is, find a symmetric matrix that has nonnegative diagonal elements but is not nonnegative definite.
- 1.39** Let A be an $m \times m$ nonnegative definite matrix, while B is an $n \times m$ matrix. Show that BAB' is a nonnegative definite matrix.
- 1.40** Define A as

$$A = \begin{bmatrix} 5 & 1 \\ 1 & 4 \end{bmatrix}.$$

Find an upper triangular square root matrix of A ; that is, find a 2×2 upper triangular matrix B satisfying $BB' = A$.

- 1.41** Use the standard normal moment generating function, $m_z(t) = e^{t^2/2}$, to show that the first six moments of the standard normal distribution are 0, 1, 0, 3, 0, and 15.

- 1.42** Use properties of expectation to show that for random variables x_1 and x_2 , and scalars α_1 , α_2 , β_1 , and β_2 ,

$$\text{cov}(\alpha_1 + \beta_1 x_1, \alpha_2 + \beta_2 x_2) = \beta_1 \beta_2 \text{cov}(x_1, x_2).$$

- 1.43** Let S be the sample covariance matrix computed from the sample $\mathbf{x}_1, \dots, \mathbf{x}_n$, where each \mathbf{x}_i is $m \times 1$. Define the $m \times n$ matrix X to be $X = (\mathbf{x}_1, \dots, \mathbf{x}_n)$. Find a matrix expression for the symmetric matrix A satisfying $S = (n-1)^{-1} X A X'$.

- 1.44** Show that if $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite, then $(\mathbf{x} - \boldsymbol{\mu})' \Omega^{-1} (\mathbf{x} - \boldsymbol{\mu}) \sim \chi_m^2$.

- 1.45** Suppose $\mathbf{x} \sim N_3(\boldsymbol{\mu}, \Omega)$, where

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix}, \quad \Omega = \begin{bmatrix} 2 & 1 & -1 \\ 1 & 2 & 1 \\ -1 & 1 & 3 \end{bmatrix},$$

and let the 3×3 matrix A and 2×3 matrix B be given by

$$A = \begin{bmatrix} 2 & 2 & 1 \\ 1 & 0 & -1 \\ 0 & 1 & -1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 1 & 0 \end{bmatrix}.$$

- (a) Find the correlation matrix of \mathbf{x} .
- (b) Determine the distribution of $u = \mathbf{1}'_3 \mathbf{x}$.
- (c) Determine the distribution of $\mathbf{v} = A\mathbf{x}$.
- (d) Determine the distribution of

$$\mathbf{w} = \begin{bmatrix} A\mathbf{x} \\ B\mathbf{x} \end{bmatrix}.$$

- (e) Which, if any, of the distributions obtained in (b), (c), and (d) are singular distributions?

- 1.46** Suppose \mathbf{x} is an $m \times 1$ random vector with mean vector $\boldsymbol{\mu}$ and covariance matrix Ω . If A is an $n \times m$ matrix of constants and \mathbf{c} is an $m \times 1$ vector of constants, give expressions for

- (a) $E[A(\mathbf{x} + \mathbf{c})]$,
- (b) $\text{var}[A(\mathbf{x} + \mathbf{c})]$.

- 1.47** Let x_1, \dots, x_m be a random sample from a normal population with mean μ and variance σ^2 , so that $\mathbf{x} = (x_1, \dots, x_m)' \sim N_m(\mu \mathbf{1}_m, \sigma^2 I_m)$.

- (a) What is the distribution of $\mathbf{u} = H\mathbf{x}$, where H is the Helmert matrix?
- (b) Show that $\sum_{i=1}^m (x_i - \bar{x})^2 = \sum_{i=2}^m u_i^2$, and use this to establish that s^2 is distributed independently of \bar{x} .

- 1.48** Use the stochastic representation given in Section 1.13 for a random vector \mathbf{x} having a contaminated normal distribution to show that $E(\mathbf{x}) = \mathbf{0}$ and $\text{var}(\mathbf{x}) = \{1 + p(\sigma^2 - 1)\}I_m$.
- 1.49** Show that if \mathbf{x} has the multivariate t distribution with n degrees of freedom as given in Section 1.13, then $E(\mathbf{x}) = \mathbf{0}$ and $\text{var}(\mathbf{x}) = \frac{n}{n-2}I_m$ if $n > 2$.

2

VECTOR SPACES

2.1 INTRODUCTION

In statistics, observations typically take the form of vectors of values of different variables; for example, for each subject in a sample, one might record height, weight, age, and so on. In estimation and hypotheses testing situations, we are usually interested in inferences regarding a vector of parameters. As a result, the topic of this chapter, vector spaces, has important applications in statistics. In addition, the concept of linearly independent and dependent vectors, which we discuss in Section 3, is very useful in the understanding and determination of the rank of a matrix.

2.2 DEFINITIONS

A vector space is a collection of vectors that satisfies some special properties; in particular, the collection is closed under the addition of vectors and under the multiplication of a vector by a scalar.

Definition 2.1 Let S be a collection of $m \times 1$ vectors satisfying the following:

- (a) If $\mathbf{x}_1 \in S$ and $\mathbf{x}_2 \in S$, then $\mathbf{x}_1 + \mathbf{x}_2 \in S$.
- (b) If $\mathbf{x} \in S$ and α is any real scalar, then $\alpha\mathbf{x} \in S$.

Then S is called a vector space in m -dimensional space. If S is a subset of T , which is another vector space in m -dimensional space, then S is called a vector subspace of T , which will be indicated by writing $S \subseteq T$.

The choice of $\alpha = 0$ in Definition 2.1(b) implies that the null vector $\mathbf{0} \in S$; that is, every vector space must contain the null vector. In fact, the set $S = \{\mathbf{0}\}$ consisting of the null vector only is itself a vector space. Note also that the two conditions (a) and (b) are equivalent to the one condition that says if $\mathbf{x}_1 \in S$, $\mathbf{x}_2 \in S$, and α_1 and α_2 are any real scalars, then $(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) \in S$. This can be easily generalized to more than two, say n , vectors; that is, if $\alpha_1, \dots, \alpha_n$ are real scalars and $\mathbf{x}_1, \dots, \mathbf{x}_n$ are vectors such that $\mathbf{x}_i \in S$, for all i , then for S to be a vector space, we must have

$$\sum_{i=1}^n \alpha_i \mathbf{x}_i \in S. \quad (2.1)$$

The left-hand side of (2.1) is called a linear combination of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$. Because a vector space is closed under the formation of linear combinations, vector spaces are sometimes also referred to as linear spaces.

The concept of linear spaces can be generalized to sets that contain elements that are not vectors. For example, if S is a set of $m \times n$ matrices, it is a linear space as long as $\alpha_1 X_1 + \alpha_2 X_2 \in S$ for all choices of $X_1 \in S$, $X_2 \in S$, and scalars α_1 and α_2 . We will use $R^{m \times n}$ to denote the linear space consisting of all $m \times n$ matrices with real components.

Example 2.1 Consider the sets of vectors given by

$$\begin{aligned} S_1 &= \{(a, 0, a)'\} : -\infty < a < \infty\}, \\ S_2 &= \{(a, b, a+b)'\} : -\infty < a < \infty, -\infty < b < \infty\}, \\ S_3 &= \{(a, a, a)'\} : a \geq 0\}. \end{aligned}$$

Let $\mathbf{x}_1 = (a_1, 0, a_1)'$ and $\mathbf{x}_2 = (a_2, 0, a_2)'$, where a_1 and a_2 are arbitrary scalars. Then $\mathbf{x}_1 \in S_1$, $\mathbf{x}_2 \in S_1$, and

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 = (\alpha_1 a_1 + \alpha_2 a_2, 0, \alpha_1 a_1 + \alpha_2 a_2)' \in S_1,$$

so that S_1 is a vector space. By a similar argument, we find that S_2 is also a vector space. Further, S_1 consists of all vectors of S_2 for which $b = 0$, so S_1 is a subset of S_2 , and thus S_1 is a vector subspace of S_2 . On the other hand, S_3 is not a vector space because, for example, if we take $\alpha = -1$ and $\mathbf{x} = (1, 1, 1)'$, then $\mathbf{x} \in S_3$ but

$$\alpha \mathbf{x} = -(1, 1, 1)' \notin S_3.$$

Every vector space with the exception of the vector space $\{0\}$ has infinitely many vectors. However, by using the process of forming linear combinations, a vector space can be associated with a finite set of vectors as long as each vector in the vector space can be expressed as some linear combination of the vectors in this set.

Definition 2.2 Let $\{x_1, \dots, x_n\}$ be a set of $m \times 1$ vectors in the vector space S . If each vector in S can be expressed as a linear combination of the vectors x_1, \dots, x_n , then the set $\{x_1, \dots, x_n\}$ is said to span or generate the vector space S , and $\{x_1, \dots, x_n\}$ is called a spanning set of S .

A spanning set for a vector space S is not uniquely defined unless $S = \{0\}$. Because $(a, b, a + b)' = a(1, 0, 1)' + b(0, 1, 1)'$, it is easy to see that $\{(1, 0, 1)', (0, 1, 1)'\}$ is a spanning set for the vector space S_2 defined in Example 2.1. However, any set of at least two vectors of this form will also be a spanning set as long as at least two of the vectors are nonnull vectors that are not scalar multiples of each other. For instance, $\{(1, 1, 2)', (1, -1, 0)', (2, 3, 5)'\}$ is also a spanning set for S_2 .

Suppose we select from the vector space S a set of vectors $\{x_1, \dots, x_n\}$. In general, we cannot be assured that every $x \in S$ is a linear combination of x_1, \dots, x_n , and so it is possible that the set $\{x_1, \dots, x_n\}$ does not span S . This set must, however, span a vector space, which is a subspace of S .

Theorem 2.1 Let $\{x_1, \dots, x_n\}$ be a set of $m \times 1$ vectors in the vector space S , and let W be the set of all possible linear combinations of these vectors; that is,

$$W = \left\{ x : x = \sum_{i=1}^n \alpha_i x_i, -\infty < \alpha_i < \infty \text{ for all } i \right\}.$$

Then W is a vector subspace of S .

Proof. Clearly, W is a subset of S because the vectors x_1, \dots, x_n are in S , and S is closed under the formation of linear combinations. To prove that W is a subspace of S , we must show that, for arbitrary vectors u and v in W and scalars a and b , $au + bv$ is in W . Because u and v are in W , by the definition of W , scalars c_1, \dots, c_n and d_1, \dots, d_n must exist such that

$$u = \sum_{i=1}^n c_i x_i, \quad v = \sum_{i=1}^n d_i x_i.$$

It then follows that

$$au + bv = a \left(\sum_{i=1}^n c_i x_i \right) + b \left(\sum_{i=1}^n d_i x_i \right) = \sum_{i=1}^n (ac_i + bd_i) x_i,$$

so that $au + bv$ is a linear combination of x_1, \dots, x_n and thus $au + bv \in W$. \square

The notions of the size or length of a vector or the distance between two vectors are important concepts when dealing with vector spaces. Although we are most familiar with the standard Euclidean formulas for length and distance, there are a variety of ways of defining length and distance. These measures of length and distance sometimes involve a product of vectors called an inner product.

Definition 2.3 Let S be a vector space. A function, $\langle \mathbf{x}, \mathbf{y} \rangle$, defined for all $\mathbf{x} \in S$ and $\mathbf{y} \in S$, is an inner product if for any \mathbf{x} , \mathbf{y} , and \mathbf{z} in S , and any scalar c :

- (a) $\langle \mathbf{x}, \mathbf{x} \rangle \geq 0$ with $\langle \mathbf{x}, \mathbf{x} \rangle = 0$ if and only if $\mathbf{x} = \mathbf{0}$.
- (b) $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{y}, \mathbf{x} \rangle$.
- (c) $\langle \mathbf{x} + \mathbf{y}, \mathbf{z} \rangle = \langle \mathbf{x}, \mathbf{z} \rangle + \langle \mathbf{y}, \mathbf{z} \rangle$.
- (d) $\langle c\mathbf{x}, \mathbf{y} \rangle = c\langle \mathbf{x}, \mathbf{y} \rangle$.

For a given inner product, define the $m \times m$ matrix A to have (i, j) th element $a_{ij} = \langle \mathbf{e}_{i,m}, \mathbf{e}_{j,m} \rangle$. Then, using properties (c) and (d) of Definition 2.3, we find that for any $m \times 1$ vectors \mathbf{x} and \mathbf{y} ,

$$\begin{aligned} \langle \mathbf{x}, \mathbf{y} \rangle &= \left\langle \left(\sum_{i=1}^m x_i \mathbf{e}_{i,m} \right), \mathbf{y} \right\rangle = \sum_{i=1}^m x_i \langle \mathbf{e}_{i,m}, \mathbf{y} \rangle \\ &= \sum_{i=1}^m x_i \left\langle \mathbf{e}_{i,m}, \sum_{j=1}^m y_j \mathbf{e}_{j,m} \right\rangle = \sum_{i=1}^m \sum_{j=1}^m x_i y_j \langle \mathbf{e}_{i,m}, \mathbf{e}_{j,m} \rangle \\ &= \sum_{i=1}^m \sum_{j=1}^m x_i y_j a_{ij} = \mathbf{x}' A \mathbf{y}. \end{aligned}$$

That is, every inner product defined on R^m can be expressed as a bilinear form $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}' A \mathbf{y}$, for some $m \times m$ matrix A . Due to properties (a) and (b) of Definition 2.3, the matrix A must be positive definite.

A useful result regarding inner products is given by the Cauchy–Schwarz inequality.

Theorem 2.2 If \mathbf{x} and \mathbf{y} are in the vector space S and $\langle \mathbf{x}, \mathbf{y} \rangle$ is an inner product defined on S , then

$$\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle, \quad (2.2)$$

with equality if and only if one of the vectors is a scalar multiple of the other.

Proof. The result is trivial if $\mathbf{x} = \mathbf{0}$ because it is easily shown that, in this case, $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle = 0$ so that (2.2) holds with equality and $\mathbf{x} = \alpha \mathbf{y}$, where $\alpha = 0$.

Suppose that $\mathbf{x} \neq \mathbf{0}$, and let $a = \langle \mathbf{x}, \mathbf{x} \rangle$, $b = 2\langle \mathbf{x}, \mathbf{y} \rangle$, and $c = \langle \mathbf{y}, \mathbf{y} \rangle$. Then using Definition 2.3, we find that for any scalar t ,

$$\begin{aligned} 0 &\leq \langle t\mathbf{x} + \mathbf{y}, t\mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle t^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle t + \langle \mathbf{y}, \mathbf{y} \rangle \\ &= at^2 + bt + c. \end{aligned} \quad (2.3)$$

Consequently, the polynomial $at^2 + bt + c$ either has a repeated real root or no real roots. This means that the discriminant $b^2 - 4ac$ must be nonpositive, and this leads to the inequality

$$b^2 \leq 4ac,$$

which simplifies to $\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle$ as is required. Now if one of the vectors is a scalar multiple of the other, then $\mathbf{y} = \alpha \mathbf{x}$ must hold for some α , and clearly this leads to equality in (2.2). Conversely, equality in (2.2) corresponds to equality in (2.3) for some t . This can only happen if $t\mathbf{x} + \mathbf{y} = \mathbf{0}$, in which case \mathbf{y} is a scalar multiple of \mathbf{x} , so the proof is complete. \square

The most common inner product is the Euclidean inner product given by $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}'\mathbf{y}$. Applying the Cauchy–Schwarz inequality to this inner product, we find that

$$\left(\sum_{i=1}^m x_i y_i \right)^2 \leq \left(\sum_{i=1}^m x_i^2 \right) \left(\sum_{i=1}^m y_i^2 \right)$$

holds for any $m \times 1$ vectors \mathbf{x} and \mathbf{y} , with equality if and only if one of these vectors is a scalar multiple of the other.

A vector norm and a distance function provide us with the means of measuring the length of a vector and the distance between two vectors.

Definition 2.4 A function $\|\mathbf{x}\|$ is a vector norm on the vector space S if, for any vectors \mathbf{x} and \mathbf{y} in S , we have

- (a) $\|\mathbf{x}\| \geq 0$,
- (b) $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$,
- (c) $\|c\mathbf{x}\| = |c| \|\mathbf{x}\|$ for any scalar c ,
- (d) $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$.

Definition 2.5 A function $d(\mathbf{x}, \mathbf{y})$ is a distance function defined on the vector space S if for any vectors \mathbf{x} , \mathbf{y} , and \mathbf{z} in S , we have

- (a) $d(\mathbf{x}, \mathbf{y}) \geq 0$,
- (b) $d(\mathbf{x}, \mathbf{y}) = 0$ if and only if $\mathbf{x} = \mathbf{y}$,
- (c) $d(\mathbf{x}, \mathbf{y}) = d(\mathbf{y}, \mathbf{x})$,
- (d) $d(\mathbf{x}, \mathbf{z}) \leq d(\mathbf{x}, \mathbf{y}) + d(\mathbf{y}, \mathbf{z})$.

Property (d) in both Definition 2.4 and Definition 2.5 is known as the triangle inequality because it is a generalization of the familiar relationship in two-dimensional geometry. One common way of defining a vector norm and a distance function is in terms of an inner product. The reader can verify that for any inner product, $\langle \mathbf{x}, \mathbf{y} \rangle$, the functions, $\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$ and $d(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x} - \mathbf{y}, \mathbf{x} - \mathbf{y} \rangle^{1/2}$ satisfy the conditions given in Definition 2.4 and Definition 2.5.

We will use R^m to denote the vector space consisting of all $m \times 1$ vectors with real components; that is, $R^m = \{(x_1, \dots, x_m)' : -\infty < x_i < \infty, i = 1, \dots, m\}$. We usually have associated with this vector space the Euclidean distance function $d_I(\mathbf{x}, \mathbf{y}) = \|\mathbf{x} - \mathbf{y}\|_2$, where $\|\mathbf{x}\|_2$ is the Euclidean norm given by

$$\|\mathbf{x}\|_2 = (\mathbf{x}'\mathbf{x})^{1/2} = \left(\sum_{i=1}^m x_i^2 \right)^{1/2}$$

and based on the Euclidean inner product $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}'\mathbf{y}$. This distance formula is a generalization of the familiar formulas that we have for distance in two- and three-dimensional geometry. The space with this distance function is called Euclidean m -dimensional space. Whenever this text works with the vector space R^m , the associated distance will be this Euclidean distance unless stated otherwise. However, in many situations in statistics, non-Euclidean distance functions are appropriate.

Example 2.2 Suppose we wish to compute the distance between the $m \times 1$ vectors \mathbf{x} and $\boldsymbol{\mu}$, where \mathbf{x} is an observation from a distribution having mean vector $\boldsymbol{\mu}$ and covariance matrix Ω . If we want to take into account the effect of the covariance structure, then Euclidean distance would not be appropriate unless $\Omega = I_m$. For example, if $m = 2$ and $\Omega = \text{diag}(0.5, 2)$, then a large value of $(x_1 - \mu_1)^2$ would be more surprising than a similar value of $(x_2 - \mu_2)^2$ because the variance of the first component of \mathbf{x} is smaller than the variance of the second component; that is, it seems reasonable in defining distance to put more weight on $(x_1 - \mu_1)^2$ than on $(x_2 - \mu_2)^2$. A more appropriate distance function is given by

$$d_\Omega(\mathbf{x}, \boldsymbol{\mu}) = \{(\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\}^{1/2},$$

and it is called the Mahalanobis distance between \mathbf{x} and $\boldsymbol{\mu}$. This function is sometimes also referred to as the distance between \mathbf{x} and $\boldsymbol{\mu}$ in the metric of Ω and is useful in a multivariate procedure known as discriminant analysis (see McLachlan, 2005, or Huberty and Olejnik, 2006). Note that if $\Omega = I_m$, then this distance function reduces to the Euclidean distance function. For $\Omega = \text{diag}(0.5, 2)$, this distance function simplifies to

$$d_\Omega(\mathbf{x}, \boldsymbol{\mu}) = \{2(x_1 - \mu_1)^2 + 0.5(x_2 - \mu_2)^2\}^{1/2}.$$

As a second illustration, suppose that again $m = 2$, but now

$$\Omega = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}.$$

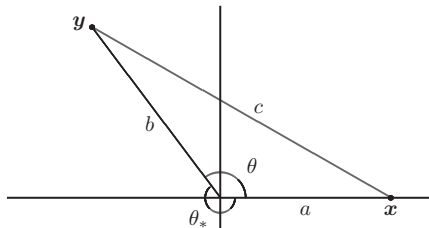


Figure 2.1 The angle between \mathbf{x} and \mathbf{y}

Because of the positive correlation, $(x_1 - \mu_1)$ and $(x_2 - \mu_2)$ will tend to have the same sign. This is reflected in the Mahalanobis distance,

$$d_{\Omega}(\mathbf{x}, \boldsymbol{\mu}) = \left(\frac{4}{3} \{ (x_1 - \mu_1)^2 + (x_2 - \mu_2)^2 - (x_1 - \mu_1)(x_2 - \mu_2) \} \right)^{1/2},$$

through the last term, which increases or decreases the distance according to whether $(x_1 - \mu_1)(x_2 - \mu_2)$ is negative or positive. In Chapter 4, we will take a closer look at the construction of this distance function.

We next consider the angle between two $m \times 1$ nonnull vectors \mathbf{x} and \mathbf{y} , where $m = 2$. We will always choose this angle to be the smaller of the two angles that can be constructed between the vectors so, for instance, in Figure 2.1, this angle is taken to be θ as opposed to $\theta_* = 2\pi - \theta$. In Figure 2.1, \mathbf{x} coincides with the first axis so $\mathbf{x} = (a, 0)'$, whereas $\mathbf{y} = (b \cos \theta, b \sin \theta)'$, where $a = (\mathbf{x}'\mathbf{x})^{1/2}$ and $b = (\mathbf{y}'\mathbf{y})^{1/2}$ are the lengths of \mathbf{x} and \mathbf{y} . The squared distance between \mathbf{x} and \mathbf{y} is then given by

$$\begin{aligned} c^2 &= (\mathbf{y} - \mathbf{x})'(\mathbf{y} - \mathbf{x}) = (b \cos \theta - a)^2 + (b \sin \theta - 0)^2 \\ &= a^2 + b^2(\cos^2 \theta + \sin^2 \theta) - 2ab \cos \theta \\ &= a^2 + b^2 - 2ab \cos \theta. \end{aligned}$$

This identity is called the law of cosines. Solving for $\cos \theta$, we get

$$\cos \theta = \frac{a^2 + b^2 - c^2}{2ab}.$$

Substituting the expressions for a , b , and c in terms of \mathbf{x} and \mathbf{y} , we obtain

$$\cos \theta = \frac{\mathbf{x}'\mathbf{y}}{(\mathbf{x}'\mathbf{x})^{1/2}(\mathbf{y}'\mathbf{y})^{1/2}},$$

an expression that is valid for all values of m and regardless of the orientation of the vectors \mathbf{x} and \mathbf{y} .

We end this section with examples of some other commonly used vector norms. The norm $\|\mathbf{x}\|_1$, called the sum norm, is defined by

$$\|\mathbf{x}\|_1 = \sum_{i=1}^m |x_i|.$$

Both the sum norm and the Euclidean norm $\|\mathbf{x}\|_2$ are members of the family of norms given by

$$\|\mathbf{x}\|_p = \left\{ \sum_{i=1}^m |x_i|^p \right\}^{1/p},$$

where $p \geq 1$. Yet another example of a vector norm, known as the infinity norm or max norm, is given by

$$\|\mathbf{x}\|_\infty = \max_{1 \leq i \leq m} |x_i|.$$

Although we have been confining attention to real vectors, these norms also serve as norms for complex vectors. However, in this case, the absolute values appearing in the expression for $\|\mathbf{x}\|_p$ are necessary even when p is even. In particular, the Euclidean norm, valid for complex as well as real vectors, is

$$\|\mathbf{x}\|_2 = \left\{ \sum_{i=1}^m |x_i|^2 \right\}^{1/2} = (\mathbf{x}^* \mathbf{x})^{1/2}.$$

2.3 LINEAR INDEPENDENCE AND DEPENDENCE

We have seen that the formation of linear combinations of vectors is a fundamental operation of vector spaces. This operation is what establishes a link between a spanning set and its vector space. In many situations, our investigation of a vector space can be reduced simply to an investigation of a spanning set for that vector space. In this case, it will be advantageous to make the spanning set as small as possible. To do this, it is first necessary to understand the concepts of linear independence and linear dependence.

Definition 2.6 The set of $m \times 1$ vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is said to be a linearly independent set if the only solution to the equation

$$\sum_{i=1}^n \alpha_i \mathbf{x}_i = \mathbf{0}$$

is given by $\alpha_1 = \dots = \alpha_n = 0$. If there are other solutions, then the set is called a linearly dependent set.

Example 2.3 Consider the three vectors $\mathbf{x}_1 = (1, 1, 1)'$, $\mathbf{x}_2 = (1, 0, -1)'$, and $\mathbf{x}_3 = (3, 2, 1)'$. To determine whether these vectors are linearly independent, we solve the system of equations $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 + \alpha_3\mathbf{x}_3 = \mathbf{0}$ or, equivalently,

$$\alpha_1 + \alpha_2 + 3\alpha_3 = 0,$$

$$\alpha_1 + 2\alpha_3 = 0,$$

$$\alpha_1 - \alpha_2 + \alpha_3 = 0.$$

These equations yield the constraints $\alpha_2 = 0.5\alpha_1$ and $\alpha_3 = -0.5\alpha_1$. Thus, for any scalar α , a solution will be given by $\alpha_1 = \alpha$, $\alpha_2 = 0.5\alpha$, and $\alpha_3 = -0.5\alpha$, and so the vectors are linearly dependent. On the other hand, any pair of these vectors are linearly independent; that is, $\{\mathbf{x}_1, \mathbf{x}_2\}$, $\{\mathbf{x}_1, \mathbf{x}_3\}$, and $\{\mathbf{x}_2, \mathbf{x}_3\}$ is each a linearly independent set of vectors.

The proof of Theorem 2.3 is left to the reader.

Theorem 2.3 Let $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ be a set of $m \times 1$ vectors. Then the following statements hold:

- (a) The set is linearly dependent if the null vector $\mathbf{0}$ is in the set.
- (b) If this set of vectors is linearly independent, any nonempty subset of it is also linearly independent.
- (c) If this set of vectors is linearly dependent, any other set containing this set as a subset is also linearly dependent.

Note that in Definition 2.6, if $n = 1$, that is, only one vector is in the set, then the set is linearly independent unless that vector is $\mathbf{0}$. If $n = 2$, the set is linearly independent unless one of the vectors is the null vector, or each vector is a nonzero scalar multiple of the other vector; that is, a set of two vectors is linearly dependent if and only if at least one of the vectors is a scalar multiple of the other. In general, we have the following result.

Theorem 2.4 The set of $m \times 1$ vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$, where $n > 1$, is a linearly dependent set if and only if at least one vector in the set can be expressed as a linear combination of the remaining vectors.

Proof. The result is obvious if one of the vectors in the set is the null vector because then the set must be linearly dependent and the $m \times 1$ null vector is a linear combination of any set of $m \times 1$ vectors. Now assume the set does not include the null vector. First suppose one of the vectors, say \mathbf{x}_n , can be expressed as a linear combination of the others; that is, we can find scalars $\alpha_1, \dots, \alpha_{n-1}$ such that $\mathbf{x}_n = \alpha_1\mathbf{x}_1 + \dots + \alpha_{n-1}\mathbf{x}_{n-1}$. But this implies that

$$\sum_{i=1}^n \alpha_i \mathbf{x}_i = \mathbf{0}, \quad (2.4)$$

if we define $\alpha_n = -1$, so the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly dependent. Conversely, now suppose that the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ are linearly dependent so that (2.4) holds for some choice of $\alpha_1, \dots, \alpha_n$ with at least one of the α_i 's, say α_n , not equal to zero. Thus, we can solve (2.4) for \mathbf{x}_n , in which case, we get

$$\mathbf{x}_n = \sum_{i=1}^{n-1} \left(\frac{-\alpha_i}{\alpha_n} \right) \mathbf{x}_i,$$

so that \mathbf{x}_n is a linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_{n-1}$. This completes the proof. \square

We end this section by proving two additional results that we will need later. Note that the first of these theorems, although stated in terms of the columns of a matrix, applies as well to the rows of a matrix.

Theorem 2.5 Consider the $m \times m$ matrix X with columns $\mathbf{x}_1, \dots, \mathbf{x}_m$. Then $|X| \neq 0$ if and only if the vectors $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly independent.

Proof. If $|X| = 0$, then $\text{rank}(X) = r < m$, and so it follows from Theorem 1.11 that nonsingular $m \times m$ matrices U and $V = [V_1 \ V_2]$ exist, with V_1 $m \times r$, such that

$$XU = V \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix} = [V_1 \ (0)].$$

But then the last column of U will give coefficients for a linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_m$, which equals the null vector. Thus, if these vectors are to be linearly independent, we must have $|X| \neq 0$. Conversely, if $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly dependent, we can find a vector $\mathbf{u} \neq \mathbf{0}$ satisfying $X\mathbf{u} = \mathbf{0}$ and then construct a nonsingular matrix U with \mathbf{u} as its last column. In this case, $XU = [W \ \mathbf{0}]$, where W is an $m \times (m-1)$ matrix and, because U is nonsingular,

$$\text{rank}(X) = \text{rank}(XU) = \text{rank}([W \ \mathbf{0}]) \leq m-1.$$

Consequently, if $|X| \neq 0$, so that $\text{rank}(X) = m$, then $\mathbf{x}_1, \dots, \mathbf{x}_m$ must be linearly independent. \square

Theorem 2.6 The set $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ of $m \times 1$ vectors is linearly dependent if $n > m$.

Proof. Consider the subset of vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_m\}$. If this is a linearly dependent set, then it follows from Theorem 2.3(c) that so is the set $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$. Thus, the proof will be complete if we can show that when $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly independent, then one of the other vectors, say \mathbf{x}_{m+1} , can be expressed as a linear combination of $\mathbf{x}_1, \dots, \mathbf{x}_m$. When $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly independent, it follows from the previous theorem that if we define X as the $m \times m$ matrix with $\mathbf{x}_1, \dots, \mathbf{x}_m$ as its columns, then $|X| \neq 0$ and so X^{-1} exists. Let $\boldsymbol{\alpha} = X^{-1}\mathbf{x}_{m+1}$ and note that $\boldsymbol{\alpha} \neq \mathbf{0}$ unless

$\mathbf{x}_{m+1} = \mathbf{0}$, in which case, the theorem is trivially true because of Theorem 2.3(a). Thus, we have

$$\sum_{i=1}^m \alpha_i \mathbf{x}_i = X\boldsymbol{\alpha} = XX^{-1}\mathbf{x}_{m+1} = \mathbf{x}_{m+1},$$

and so the set $\{\mathbf{x}_1, \dots, \mathbf{x}_{m+1}\}$ and hence also the set $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is linearly dependent. \square

2.4 MATRIX RANK AND LINEAR INDEPENDENCE

We have seen that we often work with a vector space through one of its spanning sets. In many situations, our vector space has, as a spanning set, vectors that are either the columns or rows of some matrix. In Definition 2.7, we define the terminology appropriate for such situations.

Definition 2.7 Let X be an $m \times n$ matrix. The subspace of R^n spanned by the m row vectors of X is called the row space of X . The subspace of R^m spanned by the n column vectors of X is called the column space of X .

The column space of X is sometimes also referred to as the range of X , and we will identify it by $R(X)$; that is, $R(X)$ is the vector space given by

$$R(X) = \{\mathbf{y} : \mathbf{y} = X\mathbf{a}, \mathbf{a} \in R^n\}.$$

Note that the row space of X may be written as $R(X')$.

A consequence of Theorem 2.5 is that the number of linearly independent column vectors in a matrix is identical to the rank of that matrix when it is nonsingular. Theorem 2.7 shows that this connection between the number of linearly independent columns of a matrix and the rank of that matrix always holds.

Theorem 2.7 Let X be an $m \times n$ matrix. If r is the number of linearly independent rows of X and c is the number of linearly independent columns of X , then $\text{rank}(X) = r = c$.

Proof. We will only need to prove that $\text{rank}(X) = r$ because this proof can be repeated on X' to prove that $\text{rank}(X) = c$. We will assume that the first r rows of X are linearly independent because, if they are not, elementary row transformations on X will produce such a matrix having the same rank as X . It then follows that the remaining rows of X can be expressed as linear combinations of the first r rows; that is, if X_1 is the $r \times n$ matrix consisting of the first r rows of X , then some $(m - r) \times r$ matrix A exists, such that

$$X = \begin{bmatrix} X_1 \\ AX_1 \end{bmatrix} = \begin{bmatrix} I_r \\ A \end{bmatrix} X_1.$$

Now from Theorem 2.6 we know that there can be at most r linearly independent columns in X_1 because these are $r \times 1$ vectors. Thus, we may assume that the last $n - r$ columns of X_1 can be expressed as linear combinations of the first r columns because, if this is not the case, elementary column transformations on X_1 will produce such a matrix having the same rank as X_1 . Consequently, if X_{11} is the $r \times r$ matrix with the first r columns of X_1 , then an $r \times (n - r)$ matrix B exists satisfying

$$X = \begin{bmatrix} I_r \\ A \end{bmatrix} [X_{11} \quad X_{11}B] = \begin{bmatrix} I_r \\ A \end{bmatrix} X_{11} [I_r \quad B].$$

If we define the $m \times m$ and $n \times n$ matrices U and V by

$$U = \begin{bmatrix} I_r & (0) \\ -A & I_{m-r} \end{bmatrix} \quad \text{and} \quad V = \begin{bmatrix} I_r & -B \\ (0) & I_{n-r} \end{bmatrix},$$

then we have

$$UXV = \begin{bmatrix} X_{11} & (0) \\ (0) & (0) \end{bmatrix}.$$

Because the determinant of a triangular matrix is equal to the product of its diagonal elements, we find that $|U| = |V| = 1$, so that U and V are nonsingular and thus

$$\text{rank}(X) = \text{rank}(UXV) = \text{rank}(X_{11}).$$

Finally, we must have $\text{rank}(X_{11}) = r$, because if not, by Theorem 2.5, the rows of X_{11} would be linearly dependent and this would contradict the already stated linear independence of the rows of $X_1 = [X_{11} \quad X_{11}B]$. \square

The formulation of matrix rank in terms of the number of linearly independent rows or columns of the matrix is often easier to work with than is our original definition in terms of submatrices. This is evidenced in the proof of Theorem 2.8 regarding the rank of a matrix.

Theorem 2.8 Let A be an $m \times n$ matrix. Then the following properties hold:

- (a) If B is an $n \times p$ matrix, $\text{rank}(AB) \leq \min\{\text{rank}(A), \text{rank}(B)\}$.
- (b) If B is an $m \times n$ matrix, $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$ and $\text{rank}(A + B) \geq |\text{rank}(A) - \text{rank}(B)|$.
- (c) $\text{rank}(A) = \text{rank}(A') = \text{rank}(AA') = \text{rank}(A'A)$.

Proof. Note that if B is $n \times p$, we can write

$$(AB)_{.i} = \sum_{j=1}^n b_{ji}(A)_{.j};$$

that is, each column of AB can be expressed as a linear combination of the columns of A , and so the number of linearly independent columns in AB can be no more than the number of linearly independent columns in A . Thus, $\text{rank}(AB) \leq \text{rank}(A)$. Similarly, each row of AB can be expressed as a linear combination of the rows of B from which we get $\text{rank}(AB) \leq \text{rank}(B)$, and so property (a) is proven. To prove (b), note that by using partitioned matrices, we can write

$$A + B = \begin{bmatrix} A & B \end{bmatrix} \begin{bmatrix} I_n \\ I_n \end{bmatrix}.$$

So using property (a), we find that

$$\text{rank}(A + B) \leq \text{rank}(\begin{bmatrix} A & B \end{bmatrix}) \leq \text{rank}(A) + \text{rank}(B),$$

where the final inequality follows from the easily established fact (Problem 2.26) that the number of linearly independent columns of $\begin{bmatrix} A & B \end{bmatrix}$ cannot exceed the sum of the numbers of linearly independent columns in A and in B . This establishes the first inequality in (b). Note that if we apply this inequality to A and $-B$, we also get

$$\text{rank}(A - B) \leq \text{rank}(A) + \text{rank}(B) \quad (2.5)$$

because $\text{rank}(-B) = \text{rank}(B)$. Replacing A in (2.5) by $A + B$ yields

$$\text{rank}(A + B) \geq \text{rank}(A) - \text{rank}(B),$$

and replacing B in (2.5) by $A + B$ leads to

$$\text{rank}(A + B) \geq \text{rank}(B) - \text{rank}(A).$$

Combining these two inequalities, we get the second inequality in (b). In proving (c), note that it follows immediately that $\text{rank}(A) = \text{rank}(A')$. It will suffice to prove that $\text{rank}(A) = \text{rank}(A'A)$ because this can then be used on A' to prove that $\text{rank}(A') = \text{rank}\{(A')'A'\} = \text{rank}(AA')$. If $\text{rank}(A) = r$, then a full column rank $m \times r$ matrix A_1 exists, such that, after possibly interchanging some of the columns of A , $A = \begin{bmatrix} A_1 & A_1C \end{bmatrix} = A_1 \begin{bmatrix} I_r & C \end{bmatrix}$, where C is an $r \times (n - r)$ matrix. As a result, we have

$$A'A = \begin{bmatrix} I_r \\ C' \end{bmatrix} A_1'A_1 \begin{bmatrix} I_r & C \end{bmatrix}.$$

Note that

$$EA'AE' = \begin{bmatrix} A_1'A_1 & (0) \\ (0) & (0) \end{bmatrix} \quad \text{if} \quad E = \begin{bmatrix} I_r & (0) \\ -C' & I_{n-r} \end{bmatrix},$$

and because the triangular matrix E has $|E| = 1$, E is nonsingular, so $\text{rank}(A'A) = \text{rank}(EA'AE') = \text{rank}(A_1'A_1)$. If $A_1'A_1$ is less than full rank, then by Theorem 2.5, its columns are linearly dependent, so we can find an $r \times 1$ vector $\mathbf{x} \neq \mathbf{0}$ such

that $A_1' A_1 \mathbf{x} = \mathbf{0}$, which implies that $\mathbf{x}' A_1' A_1 \mathbf{x} = (A_1 \mathbf{x})' (A_1 \mathbf{x}) = 0$. However, for any real vector \mathbf{y} , $\mathbf{y}' \mathbf{y} = 0$ only if $\mathbf{y} = \mathbf{0}$ and hence $A_1 \mathbf{x} = \mathbf{0}$. But this contradicts $\text{rank}(A_1) = r$, and so we must have $\text{rank}(A' A) = \text{rank}(A_1' A_1) = r$. \square

Theorem 2.9 gives some relationships between the rank of a partitioned matrix and the ranks of its submatrices. The proofs, which are straightforward, are left to the reader.

Theorem 2.9 Let A , B , and C be any matrices for which the partitioned matrices below are defined. Then

(a)

$$\text{rank}([A \ B]) \geq \max\{\text{rank}(A), \text{rank}(B)\},$$

(b)

$$\begin{aligned} \text{rank} \left(\begin{bmatrix} A & (0) \\ (0) & B \end{bmatrix} \right) &= \text{rank} \left(\begin{bmatrix} (0) & B \\ A & (0) \end{bmatrix} \right) \\ &= \text{rank}(A) + \text{rank}(B), \end{aligned}$$

(c)

$$\begin{aligned} \text{rank} \left(\begin{bmatrix} A & (0) \\ C & B \end{bmatrix} \right) &= \text{rank} \left(\begin{bmatrix} C & B \\ A & (0) \end{bmatrix} \right) = \text{rank} \left(\begin{bmatrix} B & C \\ (0) & A \end{bmatrix} \right) \\ &= \text{rank} \left(\begin{bmatrix} (0) & A \\ B & C \end{bmatrix} \right) \geq \text{rank}(A) + \text{rank}(B). \end{aligned}$$

Theorem 2.10 gives a useful inequality for the rank of the product of three matrices.

Theorem 2.10 Let A , B , and C be $p \times m$, $m \times n$, and $n \times q$ matrices, respectively. Then

$$\text{rank}(ABC) \geq \text{rank}(AB) + \text{rank}(BC) - \text{rank}(B).$$

Proof. It follows from Theorem 2.9(c) that

$$\text{rank} \left(\begin{bmatrix} B & BC \\ AB & (0) \end{bmatrix} \right) \geq \text{rank}(AB) + \text{rank}(BC). \quad (2.6)$$

However, because

$$\begin{bmatrix} B & BC \\ AB & (0) \end{bmatrix} = \begin{bmatrix} I_m & (0) \\ A & I_p \end{bmatrix} \begin{bmatrix} B & (0) \\ (0) & -ABC \end{bmatrix} \begin{bmatrix} I_n & C \\ (0) & I_q \end{bmatrix},$$

where, clearly, the first and last matrices on the right-hand side are nonsingular, we must also have

$$\begin{aligned} \operatorname{rank} \left(\begin{bmatrix} B & BC \\ AB & (0) \end{bmatrix} \right) &= \operatorname{rank} \left(\begin{bmatrix} B & (0) \\ (0) & -ABC \end{bmatrix} \right) \\ &= \operatorname{rank}(B) + \operatorname{rank}(ABC). \end{aligned} \quad (2.7)$$

Combining (2.6) and (2.7) we obtain the desired result. \square

A special case of Theorem 2.10 is obtained when $n = m$ and B is the $m \times m$ identity matrix. The resulting inequality gives a lower bound for the rank of a matrix product complementing the upper bound given in Theorem 2.8(a).

Corollary 2.10.1 If A is an $m \times n$ matrix and B is an $n \times p$ matrix, then

$$\operatorname{rank}(AB) \geq \operatorname{rank}(A) + \operatorname{rank}(B) - n.$$

2.5 BASES AND DIMENSION

The concept of dimension is a familiar one from geometry. For example, we recognize a line as a one-dimensional region and a plane as a two-dimensional region. In this section, we generalize this notion to any vector space. The dimension of a vector space can be determined by looking at spanning sets for that vector space. In particular, we need to be able to find the minimum number of vectors necessary for a spanning set.

Definition 2.8 Let $\{x_1, \dots, x_n\}$ be a set of $m \times 1$ vectors in a vector space S . Then this set is called a *basis* of S if it spans the vector space S and the vectors x_1, \dots, x_n are linearly independent.

Every vector space, except the vector space consisting only of the null vector 0 , has a basis. Although a basis for a vector space is not uniquely defined, a consequence of our next theorem is that the number of vectors in a basis is unique, and this is what gives us the dimension of a vector space.

Theorem 2.11 Suppose $\{x_1, \dots, x_n\}$ is a basis for the vector space S . Then

- (a) any set of more than n vectors in S must be linearly dependent,
- (b) any set of fewer than n vectors in S does not span S .

Proof. Let $\{y_1, \dots, y_k\}$ be a set of vectors in S with $k > n$. Since $\{x_1, \dots, x_n\}$ spans S and each y_i is in S , there exists an $n \times k$ matrix A such that $Y = XA$, where $X = (x_1, \dots, x_n)$ and $Y = (y_1, \dots, y_k)$. It follows from Theorem 2.8(a) that

$\text{rank}(Y) \leq \text{rank}(X) = n < k$, and this implies that the columns of Y , that is, the vectors $\mathbf{y}_1, \dots, \mathbf{y}_k$, are linearly dependent. Now assume that the set $\{\mathbf{z}_1, \dots, \mathbf{z}_p\}$ spans S . Then since each \mathbf{x}_i is in S , there exists an $p \times n$ matrix B such that $X = ZB$, where $Z = (\mathbf{z}_1, \dots, \mathbf{z}_p)$. Another application of Theorem 2.8(a) yields $n = \text{rank}(X) \leq \text{rank}(Z) \leq p$, and this establishes (b). \square

Definition 2.9 If the vector space S is $\{\mathbf{0}\}$, then the dimension of S , denoted by $\dim(S)$, is defined to be zero. Otherwise, the dimension of the vector space S is the number of vectors in any basis for S .

Example 2.4 Consider the set of $m \times 1$ vectors $\{\mathbf{e}_1, \dots, \mathbf{e}_m\}$, where for each i , \mathbf{e}_i is defined to be the vector whose only nonzero component is the i th component, which is one. Now, the linear combination of the \mathbf{e}_i 's,

$$\sum_{i=1}^m \alpha_i \mathbf{e}_i = (\alpha_1, \dots, \alpha_m)',$$

will equal $\mathbf{0}$ only if $\alpha_1 = \dots = \alpha_m = 0$, so the vectors $\mathbf{e}_1, \dots, \mathbf{e}_m$ are linearly independent. Also, if $\mathbf{x} = (x_1, \dots, x_m)'$ is an arbitrary vector in R^m , then

$$\mathbf{x} = \sum_{i=1}^m x_i \mathbf{e}_i,$$

so that $\{\mathbf{e}_1, \dots, \mathbf{e}_m\}$ spans R^m . Thus, $\{\mathbf{e}_1, \dots, \mathbf{e}_m\}$ is a basis for the m -dimensional space R^m , and in fact, any linearly independent set of m $m \times 1$ vectors will be a basis for R^m . For instance, if the $m \times 1$ vector γ_i has its first i components equal to one while the rest are all zero, then $\{\gamma_1, \dots, \gamma_m\}$ is also a basis of R^m .

Example 2.5 An implication of Theorem 2.7 is that the dimension of the column space of a matrix is the same as the dimension of the row space. However, this does not mean that the two vector spaces are the same. As a simple example, consider the matrix

$$X = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix},$$

which has rank 2. The column space of X is the two-dimensional subspace of R^3 composed of all vectors of the form $(a, b, 0)'$, whereas the row space of X is the two-dimensional subspace of R^3 containing all vectors of the form $(0, a, b)'$. If X is not square, then the column space and row space will be subspaces of different Euclidean spaces. For instance, if

$$X = \begin{bmatrix} 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 1 \end{bmatrix},$$

then the column space is R^3 , whereas the row space is the three-dimensional subspace of R^4 consisting of all vectors of the form $(a, b, c, a + b + c)'$.

If we have a spanning set for a vector space S and that spanning set is not a basis, then we could eliminate at least one vector from the set so that the reduced set is still a spanning set for S .

Theorem 2.12 If $V = \{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ spans a vector space $S \neq \{\mathbf{0}\}$ and V is a linearly dependent set of vectors, then there is a subset of V that also spans S .

Proof. Because V is a linearly dependent set, it follows from Theorem 2.4 that one of the vectors can be expressed as a linear combination of the remaining vectors. For notational convenience, we will assume that we have labeled the \mathbf{x}_i 's so that \mathbf{x}_r is such a vector; that is, scalars $\alpha_1, \dots, \alpha_{r-1}$ exist, such that $\mathbf{x}_r = \sum_{i=1}^{r-1} \alpha_i \mathbf{x}_i$. Now because V spans S , if $\mathbf{x} \in S$, scalars β_1, \dots, β_r exist, so that

$$\begin{aligned} \mathbf{x} &= \sum_{i=1}^r \beta_i \mathbf{x}_i = \beta_r \mathbf{x}_r + \sum_{i=1}^{r-1} \beta_i \mathbf{x}_i \\ &= \beta_r \sum_{i=1}^{r-1} \alpha_i \mathbf{x}_i + \sum_{i=1}^{r-1} \beta_i \mathbf{x}_i = \sum_{i=1}^{r-1} (\beta_r \alpha_i + \beta_i) \mathbf{x}_i. \end{aligned}$$

Thus, $\{\mathbf{x}_1, \dots, \mathbf{x}_{r-1}\}$ spans S and so the proof is complete. \square

Example 2.6 Consider the vector space S spanned by the vectors $\mathbf{x}_1 = (1, 1, 1)'$, $\mathbf{x}_2 = (1, 0, -1)'$, and $\mathbf{x}_3 = (3, 2, 1)'$. We saw in Example 2.3 that $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ is a linearly dependent set of vectors so that this set is not a basis for S . Because $\mathbf{x}_3 = 2\mathbf{x}_1 + \mathbf{x}_2$, we can eliminate \mathbf{x}_3 from the spanning set; that is, $\{\mathbf{x}_1, \mathbf{x}_2\}$ and $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ must span the same vector space. The set $\{\mathbf{x}_1, \mathbf{x}_2\}$ is linearly independent, and so $\{\mathbf{x}_1, \mathbf{x}_2\}$ is a basis for S and S is a two-dimensional subspace, that is, a plane in R^3 . Note that \mathbf{x}_1 is a linear combination of \mathbf{x}_2 and \mathbf{x}_3 , and \mathbf{x}_2 is a linear combination of \mathbf{x}_1 and \mathbf{x}_3 , so in this case it does not matter which \mathbf{x}_i is eliminated from the original spanning set; that is, $\{\mathbf{x}_1, \mathbf{x}_3\}$ and $\{\mathbf{x}_2, \mathbf{x}_3\}$ are also bases for S . As a related example, consider the vector space S_* spanned by the vectors $\mathbf{y}_1 = (1, 1, 1, 0)'$, $\mathbf{y}_2 = (1, 0, -1, 0)'$, $\mathbf{y}_3 = (3, 2, 1, 0)'$, and $\mathbf{y}_4 = (0, 0, 0, 1)'$. The set consisting of these four \mathbf{y}_i vectors is linearly dependent, so we will be able to eliminate one of the vectors. However, although \mathbf{y}_1 , \mathbf{y}_2 , and \mathbf{y}_3 each can be written as a linear combination of the other \mathbf{y}_i 's, \mathbf{y}_4 cannot. Thus, we can eliminate any one of the \mathbf{y}_i 's from the set except for \mathbf{y}_4 .

Every vector \mathbf{x} in a vector space can be expressed as a linear combination of the vectors in a spanning set. However, in general, more than one linear combination may yield a particular \mathbf{x} . Our next result indicates that this is not the case when the spanning set is a basis.

Theorem 2.13 Suppose the set of $m \times 1$ vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a basis for the vector space S . Then any vector $\mathbf{x} \in S$ has a unique representation as a linear combination of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$.

Proof. Because the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ span S and $\mathbf{x} \in S$, scalars $\alpha_1, \dots, \alpha_n$ must exist, such that

$$\mathbf{x} = \sum_{i=1}^n \alpha_i \mathbf{x}_i.$$

Thus, we only need to prove that the representation above is unique. Suppose it is not unique so that another set of scalars β_1, \dots, β_n exists for which

$$\mathbf{x} = \sum_{i=1}^n \beta_i \mathbf{x}_i.$$

But this result then implies that

$$\sum_{i=1}^n (\alpha_i - \beta_i) \mathbf{x}_i = \sum_{i=1}^n \alpha_i \mathbf{x}_i - \sum_{i=1}^n \beta_i \mathbf{x}_i = \mathbf{x} - \mathbf{x} = \mathbf{0}.$$

Because $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a basis, the vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ must be linearly independent and so we must have $\alpha_i - \beta_i = 0$ for all i . Thus, we must have $\alpha_i = \beta_i$, for $i = 1, \dots, n$ and so the representation is unique. \square

Some additional useful results regarding vector spaces and their bases are summarized in Theorem 2.14. The proofs are left to the reader.

Theorem 2.14 For any vector space S , the following properties hold:

- (a) If $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a set of linearly independent vectors in a vector space S and the dimension of S is n , then $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a basis for S .
- (b) If the set $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ spans the vector space S and the dimension of S is n , then the set $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ must be linearly independent and thus a basis for S .
- (c) If the vector space S has dimension n and the set of linearly independent vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ is in S , where $r < n$, then there are bases for S that contain this set as a subset.

Our final theorem in this section gives some intuitively appealing results regarding a vector space which is a subspace of another vector space.

Theorem 2.15 Suppose that S and T are subspaces of R^m with $S \subset T$. Then

- (a) $\dim(S) \leq \dim(T)$,
- (b) $S = T$ if $\dim(S) = \dim(T)$.

Proof. Let $n = \dim(S)$, $p = \dim(T)$, and suppose that $n > p$. Then a basis for S would contain n linearly independent vectors. But since $S \subset T$, these vectors are also in T leading to a contradiction of Theorem 2.11(a). Thus, we must have $n \leq p$. Now if $n = \dim(S) = \dim(T)$, then any basis for S contains exactly n vectors. But $S \subset T$, so these n linearly independent vectors are also in T and would have to form a basis for T since $\dim(T) = n$. This ensures that $S = T$. \square

2.6 ORTHONORMAL BASES AND PROJECTIONS

If each vector in a basis for a vector space S is orthogonal to every other vector in that basis, then the basis is called an orthogonal basis. In this case, the vectors can be viewed as a set of coordinate axes for the vector space S . We will find it useful also to have each vector in our basis scaled to unit length, in which case, we would have an orthonormal basis.

Suppose the set $\{x_1, \dots, x_r\}$ forms a basis for the vector space S , and we wish to obtain an orthonormal basis for S . Unless $r = 1$, an orthonormal basis is not unique so that there are many different orthonormal bases that we can construct. One method of obtaining an orthonormal basis from a given basis $\{x_1, \dots, x_r\}$ is called Gram–Schmidt orthonormalization. First, we construct the set $\{y_1, \dots, y_r\}$ of orthogonal vectors given by

$$\begin{aligned} y_1 &= x_1, \\ y_2 &= x_2 - \frac{x'_2 y_1}{y'_1 y_1} y_1, \\ &\vdots \\ y_r &= x_r - \frac{x'_r y_1}{y'_1 y_1} y_1 - \dots - \frac{x'_r y_{r-1}}{y'_{r-1} y_{r-1}} y_{r-1}, \end{aligned} \quad (2.8)$$

and then the set of orthonormal vectors $\{z_1, \dots, z_r\}$, where for each i ,

$$z_i = \frac{y_i}{(y'_i y_i)^{1/2}}.$$

Note that the linear independence of x_1, \dots, x_r guarantees the linear independence of y_1, \dots, y_r . Thus, we have Theorem 2.16.

Theorem 2.16 Every r -dimensional vector space, except the zero-dimensional space $\{0\}$, has an orthonormal basis.

If $\{z_1, \dots, z_r\}$ is a basis for the vector space S and $x \in S$, then from Theorem 2.13, we know that x can be uniquely expressed in the form $x = \alpha_1 z_1 + \dots + \alpha_r z_r$. When $\{z_1, \dots, z_r\}$ is an orthonormal basis, each of the scalars $\alpha_1, \dots, \alpha_r$ has a simple form; premultiplication of this equation for x by z'_i yields the identity $\alpha_i = z'_i x$.

Example 2.7 We will find an orthonormal basis for the three-dimensional vector space S , which has as a basis $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$, where

$$\mathbf{x}_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad \mathbf{x}_2 = \begin{bmatrix} 1 \\ -2 \\ 1 \\ -2 \end{bmatrix}, \quad \mathbf{x}_3 = \begin{bmatrix} 3 \\ 1 \\ 1 \\ -1 \end{bmatrix}.$$

The orthogonal \mathbf{y}_i 's are given by $\mathbf{y}_1 = (1, 1, 1, 1)'$,

$$\mathbf{y}_2 = \begin{bmatrix} 1 \\ -2 \\ 1 \\ -2 \end{bmatrix} - \frac{(-2)}{4} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3/2 \\ -3/2 \\ 3/2 \\ -3/2 \end{bmatrix},$$

and

$$\mathbf{y}_3 = \begin{bmatrix} 3 \\ 1 \\ 1 \\ -1 \end{bmatrix} - \frac{(4)}{(4)} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} - \frac{(6)}{(9)} \begin{bmatrix} 3/2 \\ -3/2 \\ 3/2 \\ -3/2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}.$$

Normalizing these vectors yields the orthonormal basis $\{\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3\}$, where

$$\mathbf{z}_1 = \begin{bmatrix} 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \end{bmatrix}, \quad \mathbf{z}_2 = \begin{bmatrix} 1/2 \\ -1/2 \\ 1/2 \\ -1/2 \end{bmatrix}, \quad \mathbf{z}_3 = \begin{bmatrix} 1/2 \\ 1/2 \\ -1/2 \\ -1/2 \end{bmatrix}.$$

Thus, for any $\mathbf{x} \in S$, $\mathbf{x} = \alpha_1 \mathbf{z}_1 + \alpha_2 \mathbf{z}_2 + \alpha_3 \mathbf{z}_3$, where $\alpha_i = \mathbf{x}' \mathbf{z}_i$. For instance, because $\mathbf{x}'_3 \mathbf{z}_1 = 2$, $\mathbf{x}'_3 \mathbf{z}_2 = 2$, $\mathbf{x}'_3 \mathbf{z}_3 = 2$, we have $\mathbf{x}_3 = 2\mathbf{z}_1 + 2\mathbf{z}_2 + 2\mathbf{z}_3$.

Now if S is a vector subspace of R^m and $\mathbf{x} \in R^m$, Theorem 2.17 indicates how the vector \mathbf{x} can be decomposed into the sum of a vector in S and another vector.

Theorem 2.17 Let $\{\mathbf{z}_1, \dots, \mathbf{z}_r\}$ be an orthonormal basis for some vector subspace, S , of R^m . Then each $\mathbf{x} \in R^m$ can be expressed uniquely as

$$\mathbf{x} = \mathbf{u} + \mathbf{v},$$

where $\mathbf{u} \in S$ and \mathbf{v} is a vector that is orthogonal to every vector in S .

Proof. It follows from Theorem 2.14(c) that we can find vectors $\mathbf{z}_{r+1}, \dots, \mathbf{z}_m$, so that the set $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ is an orthonormal basis for the m -dimensional Euclidean space R^m . It also follows from Theorem 2.13 that a unique set of scalars $\alpha_1, \dots, \alpha_m$ exists, such that

$$\mathbf{x} = \sum_{i=1}^m \alpha_i \mathbf{z}_i.$$

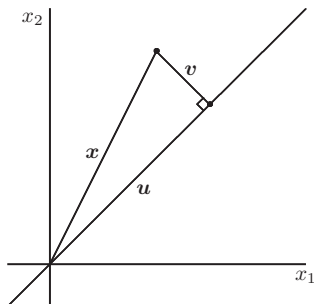


Figure 2.2 Projection of \mathbf{x} onto a one-dimensional subspace of R^2

Thus, if we let $\mathbf{u} = \alpha_1 \mathbf{z}_1 + \cdots + \alpha_r \mathbf{z}_r$ and $\mathbf{v} = \alpha_{r+1} \mathbf{z}_{r+1} + \cdots + \alpha_m \mathbf{z}_m$, we have, uniquely, $\mathbf{x} = \mathbf{u} + \mathbf{v}$, $\mathbf{u} \in S$, and \mathbf{v} will be orthogonal to every vector in S because of the orthogonality of the vectors $\mathbf{z}_1, \dots, \mathbf{z}_m$. \square

The vector \mathbf{u} in Theorem 2.17 is known as the orthogonal projection of \mathbf{x} onto S . When $m = 3$, the orthogonal projection has a simple geometrical description that allows for visualization. If, for instance, \mathbf{x} is a point in three-dimensional space and S is a two-dimensional subspace, then the orthogonal projection \mathbf{u} of \mathbf{x} will be the point of intersection of the plane S and the line that is perpendicular to S and passes through \mathbf{x} .

Example 2.8 We will look at some examples of projections in R^2 and R^3 . First consider the space S_1 spanned by the vector $\mathbf{y} = (1, 1)'$, which is a line in R^2 . Normalizing \mathbf{y} , we get $\mathbf{z}_1 = (1/\sqrt{2}, 1/\sqrt{2})'$, and the set $\{\mathbf{z}_1, \mathbf{z}_2\}$ is an orthonormal basis for R^2 , where $\mathbf{z}_2 = (1/\sqrt{2} - 1/\sqrt{2})'$. We will look at the projection of $\mathbf{x} = (1, 2)'$ onto S_1 . Solving the system of equations $\mathbf{x} = \alpha_1 \mathbf{z}_1 + \alpha_2 \mathbf{z}_2$, we find that $\alpha_1 = 3/\sqrt{2}$ and $\alpha_2 = -1/\sqrt{2}$. The orthogonal projection of \mathbf{x} onto S_1 , which is illustrated in Figure 2.2, is then given by $\mathbf{u} = \alpha_1 \mathbf{z}_1 = (3/2, 3/2)'$, whereas $\mathbf{v} = \alpha_2 \mathbf{z}_2 = (-1/2, 1/2)'$. Next consider the space S_2 spanned by the vectors $\mathbf{y}_1 = (3, 1, 1)'$ and $\mathbf{y}_2 = (1, 7, 2)'$, so that S_2 is a plane in R^3 . It is easily verified that $\{\mathbf{z}_1, \mathbf{z}_2\}$ is an orthonormal basis for S_2 , where $\mathbf{z}_1 = (3/\sqrt{11}, 1/\sqrt{11}, 1/\sqrt{11})'$ and $\mathbf{z}_2 = (-5/\sqrt{198}, 13/\sqrt{198}, 2/\sqrt{198})'$, and the set $\{\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3\}$ is an orthonormal basis for R^3 if we define $\mathbf{z}_3 = (1/\sqrt{18}, 1/\sqrt{18}, -4/\sqrt{18})'$. We will look at the projection of $\mathbf{x} = (4, 4, 4)'$ onto S_2 . Solving the system of equations $\mathbf{x} = \alpha_1 \mathbf{z}_1 + \alpha_2 \mathbf{z}_2 + \alpha_3 \mathbf{z}_3$, we find that $\alpha_1 = 20/\sqrt{11}$, $\alpha_2 = 40/\sqrt{198}$, and $\alpha_3 = -8/\sqrt{18}$. The orthogonal projection of \mathbf{x} onto S_2 , as depicted in Figure 2.3, is given by $\mathbf{u} = \alpha_1 \mathbf{z}_1 + \alpha_2 \mathbf{z}_2 = (4.44, 4.44, 2.22)'$, whereas $\mathbf{v} = (-0.44, -0.44, 1.78)'$.

The importance of the orthogonal projection \mathbf{u} in many applications arises out of the fact that it is the closest point in S to \mathbf{x} . That is, if \mathbf{y} is any other point in S and d_I is the Euclidean distance function, then $d_I(\mathbf{x}, \mathbf{u}) \leq d_I(\mathbf{x}, \mathbf{y})$. This is fairly

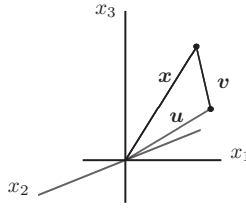


Figure 2.3 Projection of \mathbf{x} onto a two-dimensional subspace of R^3

simple to verify. Since \mathbf{u} and \mathbf{y} are in S , it follows from the decomposition $\mathbf{x} = \mathbf{u} + \mathbf{v}$ that the vector $\mathbf{u} - \mathbf{y}$ is orthogonal to $\mathbf{v} = \mathbf{x} - \mathbf{u}$ and, hence, $(\mathbf{x} - \mathbf{u})'(\mathbf{u} - \mathbf{y}) = 0$. Consequently,

$$\begin{aligned}
 \{d_I(\mathbf{x}, \mathbf{y})\}^2 &= (\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y}) \\
 &= \{(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})\}'\{(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})\} \\
 &= (\mathbf{x} - \mathbf{u})'(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})'(\mathbf{u} - \mathbf{y}) + 2(\mathbf{x} - \mathbf{u})'(\mathbf{u} - \mathbf{y}) \\
 &= (\mathbf{x} - \mathbf{u})'(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})'(\mathbf{u} - \mathbf{y}) \\
 &= \{d_I(\mathbf{x}, \mathbf{u})\}^2 + \{d_I(\mathbf{u}, \mathbf{y})\}^2,
 \end{aligned}$$

from which $d_I(\mathbf{x}, \mathbf{u}) \leq d_I(\mathbf{x}, \mathbf{y})$ follows because $\{d_I(\mathbf{u}, \mathbf{y})\}^2 \geq 0$.

Example 2.9 Simple linear regression relates a response variable y to one explanatory variable x through the model

$$y = \beta_0 + \beta_1 x + \epsilon;$$

that is, if this model is correct, then observed ordered pairs (x, y) should be clustered about some line in the x, y plane. Suppose that we have N observations, (x_i, y_i) , $i = 1, \dots, N$, and we form the $N \times 1$ vector $\mathbf{y} = (y_1, \dots, y_N)'$ and the $N \times 2$ matrix

$$X = \begin{bmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{bmatrix} = [\mathbf{1}_N \quad \mathbf{x}].$$

The least squares estimator $\hat{\beta}$ of $\beta = (\beta_0, \beta_1)'$ minimizes the sum of squared errors given by

$$(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) = (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}).$$

In Chapter 9, we will see how to find $\hat{\beta}$ using differential methods. Here we will use the geometrical properties of projections to determine $\hat{\beta}$. For any choice of $\hat{\beta}$,

$\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$ gives a point in the subspace of R^N spanned by the columns of X , that is, the plane spanned by the two vectors $\mathbf{1}_N$ and \mathbf{x} . Thus, the point $\hat{\mathbf{y}}$ that minimizes the distance from \mathbf{y} will be given by the orthogonal projection of \mathbf{y} onto this plane spanned by $\mathbf{1}_N$ and \mathbf{x} , which means that $\mathbf{y} - \hat{\mathbf{y}}$ must be orthogonal to both $\mathbf{1}_N$ and \mathbf{x} . This result leads to the two normal equations

$$\begin{aligned} 0 &= (\mathbf{y} - \hat{\mathbf{y}})' \mathbf{1}_N = \mathbf{y}' \mathbf{1}_N - \hat{\boldsymbol{\beta}}' X' \mathbf{1}_N \\ &= \sum_{i=1}^N y_i - \hat{\beta}_0 N - \hat{\beta}_1 \sum_{i=1}^N x_i, \\ 0 &= (\mathbf{y} - \hat{\mathbf{y}})' \mathbf{x} = \mathbf{y}' \mathbf{x} - \hat{\boldsymbol{\beta}}' X' \mathbf{x} \\ &= \sum_{i=1}^N x_i y_i - \hat{\beta}_0 \sum_{i=1}^N x_i - \hat{\beta}_1 \sum_{i=1}^N x_i^2, \end{aligned}$$

which when solved simultaneously for $\hat{\beta}_0$ and $\hat{\beta}_1$, yields

$$\hat{\beta}_1 = \frac{\sum_{i=1}^N x_i y_i - N \bar{x} \bar{y}}{\sum_{i=1}^N x_i^2 - N \bar{x}^2}, \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}.$$

If we want to test the hypothesis that $\beta_1 = 0$, we would consider the reduced model

$$y = \beta_0 + \epsilon,$$

and least squares estimation here only requires an estimate of β_0 . In this case, the vector of fitted values satisfies $\hat{\mathbf{y}} = \hat{\beta}_0 \mathbf{1}_N$, so for any choice of $\hat{\beta}_0$, $\hat{\mathbf{y}}$ will be given by a point on the line passing through the origin and $\mathbf{1}_N$. Thus, if $\hat{\mathbf{y}}$ is to minimize the sum of squared errors and hence the distance from \mathbf{y} , then it must be given by the orthogonal projection of \mathbf{y} onto this line. Consequently, we must have

$$0 = (\mathbf{y} - \hat{\mathbf{y}})' \mathbf{1}_N = (\mathbf{y} - \hat{\beta}_0 \mathbf{1}_N)' \mathbf{1}_N = \sum_{i=1}^N y_i - \hat{\beta}_0 N,$$

or simply

$$\hat{\beta}_0 = \bar{y}.$$

The vector \mathbf{v} in Theorem 2.17 is called the component of \mathbf{x} orthogonal to S . It is one vector belonging to what is known as the orthogonal complement of S .

Definition 2.10 Let S be a vector subspace of R^m . The orthogonal complement of S , denoted by S^\perp , is the collection of all vectors in R^m that are orthogonal to every vector in S ; that is, $S^\perp = \{\mathbf{x} : \mathbf{x} \in R^m \text{ and } \mathbf{x}' \mathbf{y} = 0 \text{ for all } \mathbf{y} \in S\}$.

Theorem 2.18 If S is a vector subspace of R^m , then its orthogonal complement S^\perp is also a vector subspace of R^m .

Proof. Suppose that $\mathbf{x}_1 \in S^\perp$ and $\mathbf{x}_2 \in S^\perp$, so that $\mathbf{x}'_1 \mathbf{y} = \mathbf{x}'_2 \mathbf{y} = 0$ for any $\mathbf{y} \in S$. Consequently, for any $\mathbf{y} \in S$ and any scalars α_1 and α_2 , we have

$$(\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2)' \mathbf{y} = \alpha_1 \mathbf{x}'_1 \mathbf{y} + \alpha_2 \mathbf{x}'_2 \mathbf{y} = 0,$$

and so $(\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2) \in S^\perp$, and thus S^\perp is a vector space. \square

A consequence of Theorem 2.19 is that if S is a vector subspace of R^m and the dimension of S is r , then the dimension of S^\perp is $m - r$.

Theorem 2.19 Suppose $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ is an orthonormal basis for R^m and $\{\mathbf{z}_1, \dots, \mathbf{z}_r\}$ is an orthonormal basis for the vector subspace S . Then $\{\mathbf{z}_{r+1}, \dots, \mathbf{z}_m\}$ is an orthonormal basis for S^\perp .

Proof. Let T be the vector space spanned by $\{\mathbf{z}_{r+1}, \dots, \mathbf{z}_m\}$. We must show that this vector space is the same as S^\perp . If $\mathbf{x} \in T$ and $\mathbf{y} \in S$, then scalars $\alpha_1, \dots, \alpha_m$ exist, such that $\mathbf{y} = \alpha_1 \mathbf{z}_1 + \dots + \alpha_r \mathbf{z}_r$ and $\mathbf{x} = \alpha_{r+1} \mathbf{z}_{r+1} + \dots + \alpha_m \mathbf{z}_m$. As a result of the orthogonality of the \mathbf{z}_i 's, $\mathbf{x}' \mathbf{y} = 0$, so $\mathbf{x} \in S^\perp$ and thus $T \subseteq S^\perp$. Conversely, suppose now that $\mathbf{x} \in S^\perp$. Because \mathbf{x} is also in R^m , scalars $\alpha_1, \dots, \alpha_m$ exist, such that $\mathbf{x} = \alpha_1 \mathbf{z}_1 + \dots + \alpha_m \mathbf{z}_m$. Now if we let $\mathbf{y} = \alpha_1 \mathbf{z}_1 + \dots + \alpha_r \mathbf{z}_r$, then $\mathbf{y} \in S$, and because $\mathbf{x} \in S^\perp$, we must have $\mathbf{x}' \mathbf{y} = \alpha_1^2 + \dots + \alpha_r^2 = 0$. But this can only happen if $\alpha_1 = \dots = \alpha_r = 0$, in which case $\mathbf{x} = \alpha_{r+1} \mathbf{z}_{r+1} + \dots + \alpha_m \mathbf{z}_m$ and so $\mathbf{x} \in T$. Thus, we also have $S^\perp \subseteq T$, and so this establishes that $T = S^\perp$. \square

2.7 PROJECTION MATRICES

The orthogonal projection of an $m \times 1$ vector \mathbf{x} onto a vector space S can be conveniently expressed in matrix form. Let $\{\mathbf{z}_1, \dots, \mathbf{z}_r\}$ be any orthonormal basis for S , whereas $\{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ is an orthonormal basis for R^m . Suppose $\alpha_1, \dots, \alpha_m$ are the constants satisfying the relationship

$$\mathbf{x} = (\alpha_1 \mathbf{z}_1 + \dots + \alpha_r \mathbf{z}_r) + (\alpha_{r+1} \mathbf{z}_{r+1} + \dots + \alpha_m \mathbf{z}_m) = \mathbf{u} + \mathbf{v},$$

where \mathbf{u} and \mathbf{v} are as previously defined. Write $\boldsymbol{\alpha} = (\alpha'_1, \alpha'_2)'$ and $Z = [Z_1 \ Z_2]$, where $\boldsymbol{\alpha}_1 = (\alpha_1, \dots, \alpha_r)'$, $\boldsymbol{\alpha}_2 = (\alpha_{r+1}, \dots, \alpha_m)'$, $Z_1 = (\mathbf{z}_1, \dots, \mathbf{z}_r)$, and $Z_2 = (\mathbf{z}_{r+1}, \dots, \mathbf{z}_m)$. Then the expression for \mathbf{x} given above can be written as

$$\mathbf{x} = Z\boldsymbol{\alpha} = Z_1\boldsymbol{\alpha}_1 + Z_2\boldsymbol{\alpha}_2;$$

that is, $\mathbf{u} = Z_1 \boldsymbol{\alpha}_1$ and $\mathbf{v} = Z_2 \boldsymbol{\alpha}_2$. As a result of the orthonormality of the \mathbf{z}_i 's, we have $Z_1' Z_1 = I_r$ and $Z_1' Z_2 = (0)$, and so

$$\begin{aligned} Z_1 Z_1' \mathbf{x} &= Z_1 Z_1' Z \boldsymbol{\alpha} = Z_1 Z_1' [Z_1 \quad Z_2] \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \end{bmatrix} \\ &= [Z_1 \quad (0)] \begin{bmatrix} \boldsymbol{\alpha}_1 \\ \boldsymbol{\alpha}_2 \end{bmatrix} = Z_1 \boldsymbol{\alpha}_1 = \mathbf{u}. \end{aligned}$$

Thus, Theorem 2.20 results.

Theorem 2.20 Suppose the columns of the $m \times r$ matrix Z_1 form an orthonormal basis for the vector space S , which is a subspace of R^m . If $\mathbf{x} \in R^m$, the orthogonal projection of \mathbf{x} onto S is given by $Z_1 Z_1' \mathbf{x}$.

The matrix $Z_1 Z_1'$ appearing in Theorem 2.20 is called the projection matrix for the vector space S and sometimes will be denoted by P_S . Similarly, $Z_2 Z_2'$ is the projection matrix for S^\perp and $ZZ' = I_m$ is the projection matrix for R^m . Since $ZZ' = Z_1 Z_1' + Z_2 Z_2'$, we have the simple equation $Z_2 Z_2' = I_m - Z_1 Z_1'$ relating the projection matrices of a vector subspace and its orthogonal complement. Although a vector space does not have a unique orthonormal basis, the projection matrix formed from these orthonormal bases is unique.

Theorem 2.21 Suppose the columns of the $m \times r$ matrices Z_1 and W_1 each form an orthonormal basis for the r -dimensional vector space S . Then $Z_1 Z_1' = W_1 W_1'$.

Proof. Each column of W_1 can be written as a linear combination of the columns of Z_1 because the columns of Z_1 span S and each column of W_1 is in S ; that is, an $r \times r$ matrix P exists, such that $W_1 = Z_1 P$. However, $Z_1' Z_1 = W_1' W_1 = I_r$, because each matrix has orthonormal columns. Thus,

$$I_r = W_1' W_1 = P' Z_1' Z_1 P = P' I_r P = P' P,$$

so that P is an orthogonal matrix. Consequently, P also satisfies $PP' = I_r$, and

$$W_1 W_1' = Z_1 P P' Z_1' = Z_1 I_r Z_1' = Z_1 Z_1',$$

so the proof is complete. \square

We will use projection matrices to take another look at the Gram–Schmidt orthonormalization procedure. The procedure takes an initial linearly independent set of vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$, which is transformed to an orthogonal set $\{\mathbf{y}_1, \dots, \mathbf{y}_r\}$, which is then transformed to an orthonormal set $\{\mathbf{z}_1, \dots, \mathbf{z}_r\}$. It is easy to verify that for $i = 1, \dots, r-1$, the vector \mathbf{y}_{i+1} can be expressed as

$$\mathbf{y}_{i+1} = \left(I_m - \sum_{j=1}^i \mathbf{z}_j \mathbf{z}_j' \right) \mathbf{x}_{i+1};$$

that is, $\mathbf{y}_{i+1} = (I_m - Z_{(i)}Z'_{(i)})\mathbf{x}_{i+1}$, where $Z_{(i)} = (\mathbf{z}_1, \dots, \mathbf{z}_i)$. Thus, the $(i+1)$ th orthogonal vector \mathbf{y}_{i+1} is obtained as the projection of the $(i+1)$ th original vector onto the orthogonal complement of the vector space spanned by the first i orthogonal vectors, $\mathbf{y}_1, \dots, \mathbf{y}_i$.

The Gram–Schmidt orthonormalization process represents one method of obtaining an orthonormal basis for a vector space S from a given basis $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$. In general, if we define the $m \times r$ matrix $X_1 = (\mathbf{x}_1, \dots, \mathbf{x}_r)$, the columns of

$$Z_1 = X_1 A \quad (2.9)$$

will form an orthonormal basis for S if A is any $r \times r$ matrix for which

$$Z'_1 Z_1 = A' X'_1 X_1 A = I_r.$$

The matrix A must be nonsingular because we must have $\text{rank}(X_1) = \text{rank}(Z_1) = r$, so A^{-1} exists, and $X'_1 X_1 = (A^{-1})' A^{-1}$ or $(X'_1 X_1)^{-1} = A A'$; that is, A is a square root matrix of $(X'_1 X_1)^{-1}$. Consequently, we can obtain an expression for the projection matrix onto the vector space S , P_S , in terms of X_1 as

$$P_S = Z_1 Z'_1 = X_1 A A' X'_1 = X_1 (X'_1 X_1)^{-1} X'_1. \quad (2.10)$$

Note that the Gram–Schmidt equations given in (2.8) can be written in matrix form as $Y_1 = X_1 T$, where $Y_1 = (\mathbf{y}_1, \dots, \mathbf{y}_r)$, $X_1 = (\mathbf{x}_1, \dots, \mathbf{x}_r)$, and T is an $r \times r$ upper triangular matrix with each diagonal element equal to 1. The normalization to produce Z_1 can then be written as $Z_1 = X_1 T D^{-1}$, where D is the diagonal matrix with the positive square root of $\mathbf{y}'_i \mathbf{y}_i$ as its i th diagonal element. Consequently, the matrix $A = T D^{-1}$ is an upper triangular matrix with positive diagonal elements. Thus, the Gram–Schmidt orthonormalization is the particular case of (2.9) in which the matrix A has been chosen to be the upper triangular square root matrix of $(X'_1 X_1)^{-1}$ having positive diagonal elements. This is commonly known as the QR factorization of A .

Example 2.10 Using the basis $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ from Example 2.7, we form the X_1 matrix

$$X_1 = \begin{bmatrix} 1 & 1 & 3 \\ 1 & -2 & 1 \\ 1 & 1 & 1 \\ 1 & -2 & -1 \end{bmatrix},$$

and it is easy to verify that

$$X'_1 X_1 = \begin{bmatrix} 4 & -2 & 4 \\ -2 & 10 & 4 \\ 4 & 4 & 12 \end{bmatrix}, \quad (X'_1 X_1)^{-1} = \frac{1}{36} \begin{bmatrix} 26 & 10 & -12 \\ 10 & 8 & -6 \\ -12 & -6 & 9 \end{bmatrix}.$$

Thus, the projection matrix for the vector space S spanned by $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ is given by

$$P_S = X_1(X_1'X_1)^{-1}X_1' = \frac{1}{4} \begin{bmatrix} 3 & 1 & 1 & -1 \\ 1 & 3 & -1 & 1 \\ 1 & -1 & 3 & 1 \\ -1 & 1 & 1 & 3 \end{bmatrix},$$

which, of course, is the same as Z_1Z_1' , where $Z_1 = (\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3)$ and $\mathbf{z}_1, \mathbf{z}_2, \mathbf{z}_3$ are the vectors obtained by the Gram–Schmidt orthonormalization in Example 2.7. Now if $\mathbf{x} = (1, 2, -1, 0)'$, then the projection of \mathbf{x} onto S is $X_1(X_1'X_1)^{-1}X_1'\mathbf{x} = \mathbf{x}$; the projection of \mathbf{x} is equal to \mathbf{x} because $\mathbf{x} = \mathbf{x}_3 - \mathbf{x}_1 - \mathbf{x}_2 \in S$. On the other hand, if $\mathbf{x} = (1, -1, 2, 1)'$, then the projection of \mathbf{x} is given by $\mathbf{u} = X_1(X_1'X_1)^{-1}X_1'\mathbf{x} = (\frac{3}{4}, -\frac{3}{4}, \frac{9}{4}, \frac{3}{4})'$. The component of \mathbf{x} orthogonal to S or, in other words, the orthogonal projection of \mathbf{x} onto S^\perp , is $\{I_4 - X_1(X_1'X_1)^{-1}X_1'\}\mathbf{x} = \mathbf{x} - X_1(X_1'X_1)^{-1}X_1'\mathbf{x} = \mathbf{x} - \mathbf{u} = (\frac{1}{4}, -\frac{1}{4}, -\frac{1}{4}, \frac{1}{4})'$, which gives us the decomposition

$$\mathbf{x} = \begin{bmatrix} 1 \\ -1 \\ 2 \\ 1 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 3 \\ -3 \\ 9 \\ 3 \end{bmatrix} + \frac{1}{4} \begin{bmatrix} 1 \\ -1 \\ -1 \\ 1 \end{bmatrix} = \mathbf{u} + \mathbf{v}$$

of Theorem 2.17.

Example 2.11 We will generalize some of the ideas of Example 2.9 to the multiple regression model

$$y = \beta_0 + \beta_1x_1 + \cdots + \beta_kx_k + \epsilon,$$

relating a response variable y to k explanatory variables, x_1, \dots, x_k . If we have N observations, this model can be written as

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where \mathbf{y} is $N \times 1$, X is $N \times (k+1)$, $\boldsymbol{\beta}$ is $(k+1) \times 1$, and $\boldsymbol{\epsilon}$ is $N \times 1$, whereas the vector of fitted values is given by

$$\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}},$$

where $\hat{\boldsymbol{\beta}}$ is an estimate of $\boldsymbol{\beta}$. Clearly, for any $\hat{\boldsymbol{\beta}}$, $\hat{\mathbf{y}}$ is a point in the subspace of R^N spanned by the columns of X . To be a least squares estimate of $\boldsymbol{\beta}$, $\hat{\boldsymbol{\beta}}$ must be such that $\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$ yields the point in this subspace closest to the vector \mathbf{y} , because this will have the sum of squared errors,

$$(\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}),$$

minimized. Thus, $X\hat{\beta}$ must be the orthogonal projection of \mathbf{y} onto the space spanned by the columns of X . If X has full column rank, then this space has projection matrix $X(X'X)^{-1}X'$, and so the required projection is

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

Premultiplying this equation by $(X'X)^{-1}X'$, we obtain the least squares estimator

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

In addition, we find that the sum of squared errors (SSE) for the fitted model $\hat{\mathbf{y}} = X\hat{\beta}$ can be written as

$$\begin{aligned} \text{SSE}_1 &= (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}) \\ &= (\mathbf{y} - X(X'X)^{-1}X'\mathbf{y})'(\mathbf{y} - X(X'X)^{-1}X'\mathbf{y}) \\ &= \mathbf{y}'(I_N - X(X'X)^{-1}X')^2\mathbf{y} \\ &= \mathbf{y}'(I_N - X(X'X)^{-1}X')\mathbf{y}, \end{aligned}$$

and so this sum of squares represents the squared length of the projection of \mathbf{y} onto the orthogonal complement of the column space of X . Suppose now that β and X are partitioned as $\beta = (\beta'_1, \beta'_2)'$ and $X = (X_1, X_2)$, where the number of columns of X_1 is the same as the number of elements in β_1 , and we wish to decide whether $\beta_2 = \mathbf{0}$. If the columns of X_1 are orthogonal to the columns of X_2 , then $X'_1X_2 = (0)$ and

$$(X'X)^{-1} = \begin{bmatrix} (X'_1X_1)^{-1} & (0) \\ (0) & (X'_2X_2)^{-1} \end{bmatrix},$$

and so $\hat{\beta}$ can be partitioned as $\hat{\beta} = (\hat{\beta}'_1, \hat{\beta}'_2)'$, where $\hat{\beta}_1 = (X'_1X_1)^{-1}X'_1\mathbf{y}$ and $\hat{\beta}_2 = (X'_2X_2)^{-1}X'_2\mathbf{y}$. Further, the sum of squared errors for the fitted model, $\hat{\mathbf{y}} = X\hat{\beta}$, can be decomposed as

$$\begin{aligned} (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}) &= \mathbf{y}'(I_N - X(X'X)^{-1}X')\mathbf{y} \\ &= \mathbf{y}'(I_N - X_1(X'_1X_1)^{-1}X'_1 - X_2(X'_2X_2)^{-1}X'_2)\mathbf{y}. \end{aligned}$$

On the other hand, the least squares estimator of β_1 in the reduced model

$$\mathbf{y} = X_1\beta_1 + \epsilon$$

is $\hat{\beta}_1 = (X'_1X_1)^{-1}X'_1\mathbf{y}$, whereas its sum of squared errors is given by

$$\begin{aligned} \text{SSE}_2 &= (\mathbf{y} - X_1\hat{\beta}_1)'(\mathbf{y} - X_1\hat{\beta}_1) \\ &= \mathbf{y}'(I_N - X_1(X'_1X_1)^{-1}X'_1)\mathbf{y}. \end{aligned}$$

Thus, the term $\text{SSE}_2 - \text{SSE}_1 = \mathbf{y}'X_2(X_2'X_2)^{-1}X_2'\mathbf{y}$ gives the reduction in the sum of squared errors attributable to the inclusion of the term $X_2\beta_2$ in the model $\mathbf{y} = X\beta + \epsilon = X_1\beta_1 + X_2\beta_2 + \epsilon$, and so its relative size will be helpful in deciding whether $\beta_2 = \mathbf{0}$. If $\beta_2 = \mathbf{0}$, then the N observations of y should be randomly clustered about the column space of X_1 in R^N with no tendency to deviate from this subspace in one direction more than in any other direction, whereas if $\beta_2 \neq \mathbf{0}$, we would expect larger deviations in directions within the column space of X_2 than in directions orthogonal to the column space of X . Now, because the dimension of the column space of X is $k + 1$, SSE_1 is the sum of squared deviations in $N - k - 1$ orthogonal directions, whereas $\text{SSE}_2 - \text{SSE}_1$ gives the sum of squared deviations in k_2 orthogonal directions, where k_2 is the number of components in β_2 . Thus, $\text{SSE}_1/(N - k - 1)$ and $(\text{SSE}_2 - \text{SSE}_1)/k_2$ should be of similar magnitudes if $\beta_2 = \mathbf{0}$, whereas the latter should be larger than the former if $\beta_2 \neq \mathbf{0}$. Consequently, a decision about β_2 can be based on the value of the statistic

$$F = \frac{(\text{SSE}_2 - \text{SSE}_1)/k_2}{\text{SSE}_1/(N - k - 1)}. \quad (2.11)$$

Using results that we will develop in Chapter 11, it can be shown that $F \sim F_{k_2, N-k-1}$ if $\epsilon \sim N_N(\mathbf{0}, \sigma^2 I_N)$ and $\beta_2 = \mathbf{0}$.

When $X_1'X_2 \neq (0)$, $(\text{SSE}_2 - \text{SSE}_1)$ does not reduce to $\mathbf{y}'X_2(X_2'X_2)^{-1}X_2'\mathbf{y}$ because, in this case, $\hat{\mathbf{y}}$ is not the sum of the projection of \mathbf{y} onto the column space of X_1 and the projection of \mathbf{y} onto the column space of X_2 . To properly assess the effect of the inclusion of the term $X_2\beta_2$ in the model, we must decompose $\hat{\mathbf{y}}$ into the sum of the projection of \mathbf{y} onto the column space of X_1 and the projection of \mathbf{y} onto the subspace of the column space of X_2 orthogonal to the column space of X_1 . This latter subspace is spanned by the columns of

$$X_{2*} = (I_N - X_1(X_1'X_1)^{-1}X_1')X_2,$$

because $(I_N - X_1(X_1'X_1)^{-1}X_1')$ is the projection matrix of the orthogonal complement of the column space of X_1 . Thus, the vector of fitted values $\hat{\mathbf{y}} = X\hat{\beta}$ can be written as

$$\hat{\mathbf{y}} = X_1(X_1'X_1)^{-1}X_1'\mathbf{y} + X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}'\mathbf{y}.$$

Further, the sum of squared errors is given by

$$\mathbf{y}'(I_N - X_1(X_1'X_1)^{-1}X_1' - X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}')\mathbf{y},$$

and the reduction in the sum of squared errors attributable to the inclusion of the term $X_2\beta_2$ in the model $\mathbf{y} = X\beta + \epsilon$ is

$$\mathbf{y}'X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}'\mathbf{y}.$$

Least squares estimators are not always unique as they have been throughout this example. For instance, let us return to the least squares estimation of β in the model

$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where now X does not have full column rank. As before, $\hat{\mathbf{y}} = X\hat{\boldsymbol{\beta}}$ will be given by the orthogonal projection of \mathbf{y} onto the space spanned by the columns of X , but the necessary projection matrix cannot be expressed as $X(X'X)^{-1}X'$, because $X'X$ is singular. If the projection matrix of the column space of X is denoted by $P_{R(X)}$, then a least squares estimator of $\boldsymbol{\beta}$ is any vector $\hat{\boldsymbol{\beta}}$ satisfying

$$X\hat{\boldsymbol{\beta}} = P_{R(X)}\mathbf{y}.$$

Since X does not have full column rank, the columns of X are linearly dependent, and so we will be able to find a nonnull vector \mathbf{a} satisfying $X\mathbf{a} = \mathbf{0}$. In this case, if $\hat{\boldsymbol{\beta}}$ is a least squares estimator of $\boldsymbol{\beta}$, so also is $\hat{\boldsymbol{\beta}} + \mathbf{a}$ because

$$X(\hat{\boldsymbol{\beta}} + \mathbf{a}) = P_{R(X)}\mathbf{y},$$

and so the least squares estimator is not unique.

We have seen that if the columns of an $m \times r$ matrix Z_1 form an orthonormal basis for a vector space S , then the projection matrix of S is given by Z_1Z_1' . Clearly this projection matrix is symmetric and, because $Z_1'Z_1 = I_r$, it is also idempotent; that is, every projection matrix is symmetric and idempotent. Theorem 2.22 proves the converse. Every symmetric idempotent matrix is a projection matrix for some vector space.

Theorem 2.22 Let P be an $m \times m$ symmetric idempotent matrix of rank r . Then an r -dimensional vector space exists that has P as its projection matrix.

Proof. From Corollary 1.11.1, an $m \times r$ matrix F and an $r \times m$ matrix G exist, such that $\text{rank}(F) = \text{rank}(G) = r$ and $P = FG$. Since P is idempotent, we have

$$FGFG = FG,$$

which implies that

$$F'FGFGG' = F'FGG'. \quad (2.12)$$

Since F and G' are full column rank, the matrices $F'F$ and GG' are nonsingular. Premultiplying (2.12) by $(F'F)^{-1}$ and postmultiplying by $(GG')^{-1}$, we obtain $GF = I_r$. Using this and the symmetry of $P = FG$, we find that

$$F = FGF = (FG)'F = G'F'F,$$

which leads to $G' = F(F'F)^{-1}$. Thus, $P = FG = F(F'F)^{-1}F'$. Comparing this with (2.10), we see that P must be the projection matrix for the vector space spanned by the columns of F . This completes the proof. \square

Example 2.12 Consider the 3×3 matrix

$$P = \frac{1}{6} \begin{bmatrix} 5 & -1 & 2 \\ -1 & 5 & 2 \\ 2 & 2 & 2 \end{bmatrix}.$$

Clearly, P is symmetric and is easily verified as idempotent, so P is a projection matrix. We will find the vector space S associated with this projection matrix. First, note that the first two columns of P are linearly independent, whereas the third column is the average of the first two columns. Thus, $\text{rank}(P) = 2$, and so the dimension of the vector space associated with P is 2. For any $\mathbf{x} \in R^3$, $P\mathbf{x}$ yields a vector in S . In particular, Pe_1 and Pe_2 are in S . These two vectors form a basis for S because they are linearly independent and the dimension of S is 2. Consequently, S contains all vectors of the form $(5a - b, 5b - a, 2a + 2b)'$.

2.8 LINEAR TRANSFORMATIONS AND SYSTEMS OF LINEAR EQUATIONS

If S is a vector subspace of R^m , with projection matrix P_S , then we have seen that for any $\mathbf{x} \in R^m$, $\mathbf{u} = \mathbf{u}(\mathbf{x}) = P_S\mathbf{x}$ is the orthogonal projection of \mathbf{x} onto S ; that is, each $\mathbf{x} \in R^m$ is transformed into a $\mathbf{u} \in S$. The function $\mathbf{u}(\mathbf{x}) = P_S\mathbf{x}$ is an example of a linear transformation of R^m into S .

Definition 2.11 Let \mathbf{u} be a function defined for all \mathbf{x} in the vector space T , such that for any $\mathbf{x} \in T$, $\mathbf{u} = \mathbf{u}(\mathbf{x}) \in S$, where S is also a vector space. Then the transformation defined by \mathbf{u} is a linear transformation of T into S if for any two scalars α_1 and α_2 and any two vectors $\mathbf{x}_1 \in T$ and $\mathbf{x}_2 \in T$,

$$\mathbf{u}(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) = \alpha_1\mathbf{u}(\mathbf{x}_1) + \alpha_2\mathbf{u}(\mathbf{x}_2).$$

We will be interested in matrix transformations of the form $\mathbf{u} = A\mathbf{x}$, where \mathbf{x} is in the subspace of R^n denoted by T , \mathbf{u} is in the subspace of R^m denoted by S , and A is an $m \times n$ matrix. This defines a transformation of T into S , and the transformation is linear because for scalars α_1 and α_2 , and $n \times 1$ vectors \mathbf{x}_1 and \mathbf{x}_2 , it follows immediately that

$$A(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) = \alpha_1A\mathbf{x}_1 + \alpha_2A\mathbf{x}_2. \quad (2.13)$$

In fact, every linear transformation can be expressed as a matrix transformation (Problem 2.40). For the orthogonal projection described at the beginning of this section, $A = P_S$, so that $n = m$, and thus, we have a linear transformation of R^m into R^m or, to be more specific, a linear transformation of R^m into S . In particular, for the multiple regression problem discussed in Example 2.11, we saw that for any

$N \times 1$ vector of observations \mathbf{y} , the vector of estimated or fitted values was given by $\hat{\mathbf{y}} = X(X'X)^{-1}X'\mathbf{y}$. Thus, because $\mathbf{y} \in R^N$ and $\hat{\mathbf{y}} \in R(X)$, we have a linear transformation of R^N into $R(X)$.

It should be obvious from (2.13) that if S is actually defined to be the set $\{\mathbf{u} : \mathbf{u} = A\mathbf{x}; \mathbf{x} \in T\}$, then T being a vector space guarantees that S will also be a vector space. In addition, if the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ span T , then the vectors $A\mathbf{x}_1, \dots, A\mathbf{x}_r$ span S . In particular, if T is R^n , then because $\mathbf{e}_1, \dots, \mathbf{e}_n$ span R^n , we find that $(A)_{\cdot 1}, \dots, (A)_{\cdot n}$ span S ; that is, S is the column space or range of A because it is spanned by the columns of A .

When the matrix A does not have full column rank, then vectors \mathbf{x} will exist, other than the null vector, which satisfy $A\mathbf{x} = \mathbf{0}$. The set of all such vectors is called the nullspace of the transformation $A\mathbf{x}$ or simply the null space of the matrix A .

Theorem 2.23 Let the linear transformation of R^n into S be given by $\mathbf{u} = A\mathbf{x}$, where $\mathbf{x} \in R^n$ and A is an $m \times n$ matrix. Then the null space of A , given by the set

$$N(A) = \{\mathbf{x} : A\mathbf{x} = \mathbf{0}, \mathbf{x} \in R^n\},$$

is a vector space.

Proof. Let \mathbf{x}_1 and \mathbf{x}_2 be in $N(A)$ so that $A\mathbf{x}_1 = A\mathbf{x}_2 = \mathbf{0}$. Then, for any scalars α_1 and α_2 , we have

$$A(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) = \alpha_1 A\mathbf{x}_1 + \alpha_2 A\mathbf{x}_2 = \alpha_1(\mathbf{0}) + \alpha_2(\mathbf{0}) = \mathbf{0},$$

so that $(\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2) \in N(A)$ and, hence, $N(A)$ is a vector space. \square

The null space of a matrix A is related to the concept of orthogonal complements discussed in Section 2.6. In fact, the null space of the matrix A is the same as the orthogonal complement of the row space of A . Similarly, the null space of the matrix A' is the same as the orthogonal complement of the column space of A . Theorem 2.24 is an immediate consequence of Theorem 2.19.

Theorem 2.24 Let A be an $m \times n$ matrix. If the dimension of the row space of A is r_1 and the dimension of the null space of A is r_2 , then $r_1 + r_2 = n$.

Since the rank of the matrix A is equal to the dimension of the row space of A , the result above can be equivalently expressed as

$$\text{rank}(A) = n - \dim\{N(A)\}. \quad (2.14)$$

This connection between the rank of a matrix and the dimension of the null space of that matrix can be useful to us in determining the rank of a matrix in certain situations.

Example 2.13 To illustrate the utility of (2.14), we will give an alternative proof of the identity $\text{rank}(A) = \text{rank}(A'A)$, which was given as Theorem 2.8(c).

Suppose \mathbf{x} is in the null space of A so that $A\mathbf{x} = \mathbf{0}$. Then, clearly, we must have $A'A\mathbf{x} = \mathbf{0}$, which implies that \mathbf{x} is also in the null space of $A'A$, so it follows that $\dim\{N(A)\} \leq \dim\{N(A'A)\}$, or equivalently,

$$\text{rank}(A) \geq \text{rank}(A'A). \quad (2.15)$$

On the other hand, if \mathbf{x} is in the null space of $A'A$, then $A'A\mathbf{x} = \mathbf{0}$. Premultiplying by \mathbf{x}' yields $\mathbf{x}'A'A\mathbf{x} = 0$, which is satisfied only if $A\mathbf{x} = \mathbf{0}$. Thus, \mathbf{x} is also in the null space of A so that $\dim\{N(A)\} \geq \dim\{N(A'A)\}$, or

$$\text{rank}(A) \leq \text{rank}(A'A). \quad (2.16)$$

Combining (2.15) and (2.16), we get $\text{rank}(A) = \text{rank}(A'A)$.

When A is an $m \times m$ nonsingular matrix and $\mathbf{x} \in R^m$, then $\mathbf{u} = A\mathbf{x}$ defines a one-to-one transformation of R^m onto R^m . One way of viewing this transformation is as the movement of each point in R^m to another point in R^m . Alternatively, we can view the transformation as a change of coordinate axes. For instance, if we start with the standard coordinate axes, which are given by the columns $\mathbf{e}_1, \dots, \mathbf{e}_m$ of the identity matrix I_m , then, because for any $\mathbf{x} \in R^m$, $\mathbf{x} = x_1\mathbf{e}_1 + \dots + x_m\mathbf{e}_m$, the components of \mathbf{x} give the coordinates of the point \mathbf{x} relative to these standard coordinate axes. On the other hand, if $\mathbf{x}_1, \dots, \mathbf{x}_m$ is another basis for R^m , then from Theorem 2.13, scalars u_1, \dots, u_m exist, so that with $\mathbf{u} = (u_1, \dots, u_m)'$ and $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, we have

$$\mathbf{x} = \sum_{i=1}^m u_i \mathbf{x}_i = X\mathbf{u};$$

that is, $\mathbf{u} = (u_1, \dots, u_m)'$ gives the coordinates of the point \mathbf{x} relative to the coordinate axes $\mathbf{x}_1, \dots, \mathbf{x}_m$. The transformation from the standard coordinate system to the one with axes $\mathbf{x}_1, \dots, \mathbf{x}_m$ is then given by the matrix transformation $\mathbf{u} = A\mathbf{x}$, where $A = X^{-1}$. Note that the squared Euclidean distance of \mathbf{u} from the origin,

$$\mathbf{u}'\mathbf{u} = (A\mathbf{x})'(A\mathbf{x}) = \mathbf{x}'A'A\mathbf{x},$$

will be the same as the squared Euclidean distance of \mathbf{x} from the origin for every choice of \mathbf{x} if and only if A , and hence, also X , is an orthogonal matrix. In this case, $\mathbf{x}_1, \dots, \mathbf{x}_m$ forms an orthonormal basis for R^m , and so the transformation has replaced the standard coordinate axes by a new set of orthogonal axes given by $\mathbf{x}_1, \dots, \mathbf{x}_m$.

Example 2.14 Orthogonal transformations are of two types according to whether the determinant of A is $+1$ or -1 . If $|A| = 1$, then the new axes can be obtained by a rotation of the standard axes. For example, for a fixed angle θ , let

$$A = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

so that $|A| = \cos^2 \theta + \sin^2 \theta = 1$. The transformation given by $\mathbf{u} = A\mathbf{x}$ transforms the standard axes $\mathbf{e}_1, \mathbf{e}_2, \mathbf{e}_3$ to the new axes $\mathbf{x}_1 = (\cos \theta, -\sin \theta, 0)'$, $\mathbf{x}_2 = (\sin \theta, \cos \theta, 0)'$, $\mathbf{x}_3 = \mathbf{e}_3$, and this simply represents a rotation of \mathbf{e}_1 and \mathbf{e}_2 through an angle of θ . If instead we have

$$A = \begin{bmatrix} \cos \theta & -\sin \theta & 0 \\ \sin \theta & \cos \theta & 0 \\ 0 & 0 & -1 \end{bmatrix},$$

then $|A| = (\cos^2 \theta + \sin^2 \theta) \cdot (-1) = -1$. Now the transformation given by $\mathbf{u} = A\mathbf{x}$ transforms the standard axes to the new axes $\mathbf{x}_1 = (\cos \theta, -\sin \theta, 0)'$, $\mathbf{x}_2 = (\sin \theta, \cos \theta, 0)'$, and $\mathbf{x}_3 = -\mathbf{e}_3$; these axes are obtained by a rotation of \mathbf{e}_1 and \mathbf{e}_2 through an angle of θ followed by a reflection of \mathbf{e}_3 about the $\mathbf{x}_1, \mathbf{x}_2$ plane.

Although orthogonal transformations are common, situations occur in which non-singular, nonorthogonal transformations are useful.

Example 2.15 Suppose we have several three-dimensional vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ that are observations from distributions, each having the same positive definite covariance matrix Ω . If we are interested in how these vectors differ from one another, then a plot of the points in R^3 may be helpful. However, as discussed in Example 2.2, if Ω is not the identity matrix, then the Euclidean distance is not appropriate, and so it becomes difficult to compare and interpret the observed differences among the r points. This difficulty can be resolved by an appropriate transformation. We will see in Chapter 3 and Chapter 4 that because Ω is positive definite, a nonsingular matrix T exists, which satisfies $\Omega = TT'$. If we let $\mathbf{u}_i = T^{-1}\mathbf{x}_i$, then the Mahalanobis distance, which was defined in Example 2.2, between \mathbf{x}_i and \mathbf{x}_j is

$$\begin{aligned} d_\Omega(\mathbf{x}_i, \mathbf{x}_j) &= \{(\mathbf{x}_i - \mathbf{x}_j)' \Omega^{-1} (\mathbf{x}_i - \mathbf{x}_j)\}^{1/2} \\ &= \{(\mathbf{x}_i - \mathbf{x}_j)' T^{-1'} T^{-1} (\mathbf{x}_i - \mathbf{x}_j)\}^{1/2} \\ &= \{(T^{-1}\mathbf{x}_i - T^{-1}\mathbf{x}_j)' (T^{-1}\mathbf{x}_i - T^{-1}\mathbf{x}_j)\}^{1/2} \\ &= \{(\mathbf{u}_i - \mathbf{u}_j)' (\mathbf{u}_i - \mathbf{u}_j)\}^{1/2} = d_I(\mathbf{u}_i, \mathbf{u}_j), \end{aligned}$$

whereas the variance of \mathbf{u}_i is given by

$$\begin{aligned} \text{var}(\mathbf{u}_i) &= \text{var}(T^{-1}\mathbf{x}_i) = T^{-1} \{\text{var}(\mathbf{x}_i)\} T^{-1'} \\ &= T^{-1} \Omega T^{-1'} = T^{-1} T T' T^{-1'} = I_3. \end{aligned}$$

That is, the transformation $\mathbf{u}_i = T^{-1}\mathbf{x}_i$ produces vectors for which the Euclidean distance function is an appropriate measure of distance between points.

In Example 2.16 and Example 2.17, we discuss some transformations that are sometimes useful in regression analysis.

Example 2.16 A simple transformation that is useful in some situations is one that centers a collection of numbers at the origin. For instance, if \bar{x} is the mean of the components of $\mathbf{x} = (x_1, \dots, x_N)'$, then the average of each component of

$$\mathbf{v} = (I_N - N^{-1}\mathbf{1}_N\mathbf{1}_N')\mathbf{x} = \begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_N - \bar{x} \end{bmatrix}$$

is 0. This transformation is sometimes used in a regression analysis to center each of the explanatory variables. Thus, the multiple regression model

$$\begin{aligned} \mathbf{y} &= X\boldsymbol{\beta} + \boldsymbol{\epsilon} = [\mathbf{1}_N \quad X_1] \begin{bmatrix} \beta_0 \\ \boldsymbol{\beta}_1 \end{bmatrix} + \boldsymbol{\epsilon} \\ &= \beta_0\mathbf{1}_N + X_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon} \end{aligned}$$

can be re-expressed as

$$\begin{aligned} \mathbf{y} &= \beta_0\mathbf{1}_N + \{N^{-1}\mathbf{1}_N\mathbf{1}_N' + (I_N - N^{-1}\mathbf{1}_N\mathbf{1}_N')\}X_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon} \\ &= \gamma_0\mathbf{1}_N + V_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon} = V\boldsymbol{\gamma} + \boldsymbol{\epsilon}, \end{aligned}$$

where $V = [\mathbf{1}_N \quad V_1] = [\mathbf{1}_N \quad (I_N - N^{-1}\mathbf{1}_N\mathbf{1}_N')X_1]$ and $\boldsymbol{\gamma} = (\gamma_0, \boldsymbol{\beta}_1)' = (\beta_0 + N^{-1}\mathbf{1}_N'X_1\boldsymbol{\beta}_1, \boldsymbol{\beta}_1)'$. Because the columns of V_1 are orthogonal to $\mathbf{1}_N$, the least squares estimator of $\boldsymbol{\gamma}$ simplifies to

$$\begin{aligned} \hat{\boldsymbol{\gamma}} &= \begin{bmatrix} \hat{\gamma}_0 \\ \hat{\boldsymbol{\beta}}_1 \end{bmatrix} = (V'V)^{-1}V'\mathbf{y} \\ &= \begin{bmatrix} N^{-1} & \mathbf{0}' \\ \mathbf{0} & (V_1'V_1)^{-1} \end{bmatrix} \begin{bmatrix} \sum y_i \\ V_1'\mathbf{y} \end{bmatrix} \\ &= \begin{bmatrix} \bar{y} \\ (V_1'V_1)^{-1}V_1'\mathbf{y} \end{bmatrix}. \end{aligned}$$

Thus, $\hat{\gamma}_0 = \bar{y}$. The estimator, $\hat{\boldsymbol{\beta}}_1$, can be conveniently expressed in terms of the sample covariance matrix computed from the N $(k+1) \times 1$ vectors that form the rows of the matrix $[\mathbf{y} \quad X_1]$. If we denote this covariance matrix by S and partition it as

$$S = \begin{bmatrix} s_{11} & \mathbf{s}_{21}' \\ \mathbf{s}_{21} & S_{22} \end{bmatrix},$$

then $(N-1)^{-1}V_1'V_1 = S_{22}$ and, because $V_1'\mathbf{1}_N = \mathbf{0}$,

$$(N-1)^{-1}V_1'\mathbf{y} = (N-1)^{-1}V_1'(\mathbf{y} - \bar{y}\mathbf{1}_N) = \mathbf{s}_{21}.$$

Consequently, $\hat{\beta}_1 = S_{22}^{-1} s_{21}$. Yet another adjustment to the original regression model involves the standardization of the explanatory variables. In this case, the model becomes

$$\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \delta_1 + \epsilon = Z \delta + \epsilon,$$

where $\delta = (\delta_0, \delta_1')'$, $Z = [\mathbf{1}_N \quad Z_1]$, $\delta_0 = \gamma_0$, $Z_1 = V_1 D_{S_{22}}^{-1/2}$, and $\delta_1 = D_{S_{22}}^{1/2} \beta_1$. The least squares estimators are $\hat{\delta}_0 = \bar{y}$ and $\hat{\delta}_1 = s_{11}^{1/2} R_{22}^{-1} r_{21}$, where we have partitioned the correlation matrix R , computed from the $N(k+1) \times 1$ vectors that form the rows of the matrix $[\mathbf{y} \quad X_1]$, in a fashion similar to that of S .

The centering of explanatory variables, discussed previously, involves a linear transformation on the columns of X_1 . In some situations, it is advantageous to employ a linear transformation on the rows of X_1 , V_1 , or Z_1 . For instance, suppose that T is a $k \times k$ nonsingular matrix, and we define $W_1 = Z_1 T$, $\alpha_0 = \delta_0$, and $\alpha_1 = T^{-1} \delta_1$, so that the model

$$\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \delta_1 + \epsilon = Z \delta + \epsilon$$

can be written as

$$\mathbf{y} = \alpha_0 \mathbf{1}_N + W_1 \alpha_1 + \epsilon = W \alpha + \epsilon,$$

where $W = [\mathbf{1}_N \quad W_1]$. This second model uses a different set of explanatory variables than the first; its i th explanatory variable is a linear combination of the explanatory variables of the first model with the coefficients given by the i th column of T . However, the two models yield equivalent results in terms of the fitted values. To see this, let

$$T_* = \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & T \end{bmatrix},$$

so that $W = Z T_*$, and note that the vector of fitted values from the second model,

$$\begin{aligned} \hat{\mathbf{y}} &= W \hat{\alpha} = W(W'W)^{-1}W'\mathbf{y} = Z T_* (T_*' Z' Z T_*)^{-1} T_*' Z' \mathbf{y} \\ &= Z T_* T_*^{-1} (Z' Z)^{-1} T_*^{-1'} T_*' Z' \mathbf{y} = Z (Z' Z)^{-1} Z' \mathbf{y}, \end{aligned}$$

is the same as that obtained from the first model.

Example 2.17 Consider the multiple regression model

$$\mathbf{y} = X\beta + \epsilon,$$

where now $\text{var}(\epsilon) \neq \sigma^2 I_N$. In this case, our previous estimator, $\hat{\beta} = (X'X)^{-1}X'\mathbf{y}$, is still the least squares estimator of β , but it does not possess certain optimality properties, one of which is illustrated later in Example 3.14, that hold when $\text{var}(\epsilon) = \sigma^2 I_N$. In this example, we will consider the situation in which the ϵ_i 's are still uncorrelated, but their variances are not all the same. Thus, $\text{var}(\epsilon) = \Omega = \sigma^2 C$, where

$C = \text{diag}(c_1^2, \dots, c_N^2)$ and the c_i 's are known constants. This special regression problem is sometimes referred to as weighted least squares regression. The weighted least squares estimator of β is obtained by making a simple transformation so that ordinary least squares regression applies to the transformed model. Define the matrix $C^{-1/2} = \text{diag}(c_1^{-1}, \dots, c_N^{-1})$ and transform the original regression problem by pre-multiplying the model equation by $C^{-1/2}$; the new model equation is

$$C^{-1/2}\mathbf{y} = C^{-1/2}X\beta + C^{-1/2}\epsilon$$

or, equivalently,

$$\mathbf{y}_* = X_*\beta + \epsilon_*,$$

where $\mathbf{y}_* = C^{-1/2}\mathbf{y}$, $X_* = C^{-1/2}X$, and $\epsilon_* = C^{-1/2}\epsilon$. The covariance matrix of ϵ_* is

$$\begin{aligned} \text{var}(\epsilon_*) &= \text{var}(C^{-1/2}\epsilon) = C^{-1/2}\text{var}(\epsilon)C^{-1/2} \\ &= C^{-1/2}\{\sigma^2 C\}C^{-1/2} = \sigma^2 I_N. \end{aligned}$$

Thus, for the transformed model, ordinary least squares regression applies, and so the least squares estimator of β can be expressed as

$$\hat{\beta} = (X_*'X_*)^{-1}X_*'\mathbf{y}_*.$$

Rewriting this equation in the original model terms X and \mathbf{y} , we get

$$\begin{aligned} \hat{\beta} &= (X'C^{-1/2}C^{-1/2}X)^{-1}X'C^{-1/2}C^{-1/2}\mathbf{y} \\ &= (X'C^{-1}X)^{-1}X'C^{-1}\mathbf{y}. \end{aligned}$$

A common application related to linear transformations is one in which the matrix A and vector \mathbf{u} consist of known constants, whereas \mathbf{x} is a vector of variables, and we wish to determine all \mathbf{x} for which $A\mathbf{x} = \mathbf{u}$; that is, we want to find the simultaneous solutions x_1, \dots, x_n to the system of m equations

$$\begin{aligned} a_{11}x_1 + \dots + a_{1n}x_n &= u_1 \\ &\vdots \\ a_{m1}x_1 + \dots + a_{mn}x_n &= u_m. \end{aligned}$$

For instance, in Example 2.11, we saw that the least squares estimator of the parameter vector β in the multiple regression model satisfies the equation, $X\hat{\beta} = X(X'X)^{-1}X'\mathbf{y}$; that is, here $A = X$, $\mathbf{u} = X(X'X)^{-1}X'\mathbf{y}$, and $\mathbf{x} = \hat{\beta}$. In general, if $\mathbf{u} = \mathbf{0}$, then this system of equations is referred to as a homogeneous system, and the set of all solutions to $A\mathbf{x} = \mathbf{u}$, in this case, is simply given by the

null space of A . Consequently, if A has full column rank, then $\mathbf{x} = \mathbf{0}$ is the only solution, whereas infinitely many solutions exist if A has less than full column rank. A nonhomogeneous system of linear equations is one that has $\mathbf{u} \neq \mathbf{0}$. Although a homogeneous system always has at least one solution, $\mathbf{x} = \mathbf{0}$, a nonhomogeneous system might not have any solutions. A system of linear equations that has no solutions is called an inconsistent system of equations, whereas a system with solutions is referred to as a consistent system. If $\mathbf{u} \neq \mathbf{0}$ and $A\mathbf{x} = \mathbf{u}$ holds for some \mathbf{x} , then \mathbf{u} must be a linear combination of the columns of A ; that is, the nonhomogeneous system of equations $A\mathbf{x} = \mathbf{u}$ is consistent if and only if \mathbf{u} is in the column space of A .

The mathematics involved in solving systems of linear equations is most conveniently handled using matrix methods. For example, consider one of the simplest nonhomogeneous systems of linear equations in which the matrix A is square and nonsingular. In this case, because A^{-1} exists, we find that the system $A\mathbf{x} = \mathbf{u}$ has a solution that is unique and is given by $\mathbf{x} = A^{-1}\mathbf{u}$. Similarly, when the matrix A is singular or not even square, matrix methods can be used to determine whether the system is consistent, and if so, the solutions can be given as matrix expressions. The results regarding the solution of a general system of linear equations will be developed in Chapter 6.

The focus in this section has been on linear transformations of vectors, but the concept is easily extended to matrices. Suppose T is a linear space of $p \times q$ matrices, and let X and Z be $m \times p$ and $q \times n$ matrices, respectively. Then for $B \in T$, $Y = XBZ$ is a linear transformation of the linear space $T \subset R^{p \times q}$ onto the linear space $S = \{Y : Y = XBZ, B \in T\} \subset R^{m \times n}$. When T is $R^{p \times q}$, this transformation yields the range $R(X, Z)$; that is, $R(X, Z) = \{Y : Y = XBZ, B \in R^{p \times q}\}$, so that this linear space consists of all matrices Y whose columns are in $R(X)$ and whose rows are in $R(Z')$. The null space of the linear transformation $Y = XBZ$ is given by $N(X, Z) = \{B : XBZ = (0), B \in R^{p \times q}\} \subset R^{p \times q}$.

Our final example of this section concerns a statistical analysis that utilizes a matrix transformation of the form $Y = XBZ$.

Example 2.18 In this example, we look at a statistical model commonly known as a growth curve model. One response variable is recorded for each of n individuals at m different points in time, t_1, \dots, t_m . The expected value of the response for any individual at time t_i is modeled as a polynomial in time of degree $p - 1$ for some choice of p ; that is, this expected value has the form $b_0 + b_1 t_i + \dots + b_{p-1} t_i^{p-1}$, where b_0, \dots, b_{p-1} are unknown parameters. The response matrix Y is $m \times n$ with (i, j) th component given by the response for the j th individual at time t_i , and the growth curve model has the matrix form $Y = XBZ + E$. The (i, j) th element of the $m \times p$ matrix X is t_i^{j-1} and the $m \times n$ matrix E is a matrix of random errors. Numerous designs can be achieved through appropriate choices of the $p \times q$ parameter matrix B and the $q \times n$ design matrix Z . For instance, if all n individuals come from a common group, then we would take $q = 1$, $B = (b_0, b_1, \dots, b_{p-1})'$, and $Z = \mathbf{1}'_n$. If, on the other hand, there are g groups and n_i of the individuals are from the i th

group, we would have $q = g$,

$$B = \begin{bmatrix} b_{01} & \cdots & b_{0g} \\ \vdots & & \vdots \\ b_{p-1,1} & \cdots & b_{p-1,g} \end{bmatrix}, \quad Z = \begin{bmatrix} \mathbf{1}'_{n_1} & \mathbf{0}' & \cdots & \mathbf{0}' \\ \mathbf{0}' & \mathbf{1}'_{n_2} & \cdots & \mathbf{0}' \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0}' & \mathbf{0}' & \cdots & \mathbf{1}'_{n_g} \end{bmatrix}.$$

The method of least squares can be used to find an estimator for B . We will show that this estimator is given by $\hat{B} = (X'X)^{-1}X'YZ'(ZZ')^{-1}$ when $\text{rank}(X) = p$ and $\text{rank}(Z) = q$. Note that the sum of squared errors corresponding to a particular B matrix is $\text{tr}\{(Y - XBZ)(Y - XBZ)'\}$ and

$$\begin{aligned} & \text{tr}\{(Y - XBZ)(Y - XBZ)'\} \\ &= \text{tr}\{[(Y - X\hat{B}Z) + (X\hat{B}Z - XBZ)]\{[(Y - X\hat{B}Z) + (X\hat{B}Z - XBZ)]'\} \\ &= \text{tr}\{(Y - X\hat{B}Z)(Y - X\hat{B}Z)'\} + \text{tr}\{(X\hat{B}Z - XBZ)(X\hat{B}Z - XBZ)'\} \\ &\quad + 2\text{tr}\{(Y - X\hat{B}Z)(X\hat{B}Z - XBZ)'\} \\ &\geq \text{tr}\{(Y - X\hat{B}Z)(Y - X\hat{B}Z)'\} + 2\text{tr}\{(Y - X\hat{B}Z)(X\hat{B}Z - XBZ)'\}, \end{aligned}$$

where the inequality follows from the fact that the trace of a nonnegative definite matrix is nonnegative. But since $X'X\hat{B}ZZ' = X'YZ'$, we have

$$\begin{aligned} \text{tr}\{(Y - X\hat{B}Z)(X\hat{B}Z - XBZ)'\} &= \text{tr}\{X'(Y - X\hat{B}Z)Z'(\hat{B} - B)'\} \\ &= \text{tr}\{(X'YZ' - X'X\hat{B}ZZ')(\hat{B} - B)'\} \\ &= \text{tr}\{(X'YZ' - X'YZ')(\hat{B} - B)'\} = 0, \end{aligned}$$

so it follows that

$$\text{tr}\{(Y - XBZ)(Y - XBZ)'\} \geq \text{tr}\{(Y - X\hat{B}Z)(Y - X\hat{B}Z)'\},$$

which confirms that \hat{B} is the least squares estimator of B .

2.9 THE INTERSECTION AND SUM OF VECTOR SPACES

In this section, we discuss some common ways of forming a vector subspace from two or more given subspaces. The first of these uses a familiar operation from set theory.

Definition 2.12 Let S_1 and S_2 be vector subspaces of R^m . The intersection of S_1 and S_2 , denoted by $S_1 \cap S_2$, is the vector subspace given as

$$S_1 \cap S_2 = \{\mathbf{x} \in R^m : \mathbf{x} \in S_1 \text{ and } \mathbf{x} \in S_2\}.$$

Note that this definition says that the set $S_1 \cap S_2$ is a vector subspace if S_1 and S_2 are vector subspaces, which follows from the fact that if \mathbf{x}_1 and \mathbf{x}_2 are in $S_1 \cap S_2$, then $\mathbf{x}_1 \in S_1$, $\mathbf{x}_2 \in S_1$ and $\mathbf{x}_1 \in S_2$, $\mathbf{x}_2 \in S_2$. Thus, because S_1 and S_2 are vector spaces, for any scalars α_1 and α_2 , $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2$ will be in S_1 and S_2 and, hence, in $S_1 \cap S_2$. Definition 2.12 can be generalized in an obvious fashion to the intersection, $S_1 \cap \cdots \cap S_r$, of the r vector spaces S_1, \dots, S_r .

A second set operation, which combines the elements of S_1 and S_2 , is the union; that is, the union of S_1 and S_2 is given by

$$S_1 \cup S_2 = \{\mathbf{x} \in R^m : \mathbf{x} \in S_1 \text{ or } \mathbf{x} \in S_2\}.$$

If S_1 and S_2 are vector subspaces, then $S_1 \cup S_2$ will also be a vector subspace only if $S_1 \subseteq S_2$ or $S_2 \subseteq S_1$. It can be easily shown that the following combination of S_1 and S_2 yields the vector space containing $S_1 \cup S_2$ with the smallest possible dimension.

Definition 2.13 If S_1 and S_2 are vector subspaces of R^m , then the sum of S_1 and S_2 , denoted by $S_1 + S_2$, is the vector space given by

$$S_1 + S_2 = \{\mathbf{x}_1 + \mathbf{x}_2 : \mathbf{x}_1 \in S_1, \mathbf{x}_2 \in S_2\}.$$

Again our definition can be generalized to $S_1 + \cdots + S_r$, the sum of the r vector spaces S_1, \dots, S_r . The proof of Theorem 2.25 has been left as an exercise.

Theorem 2.25 If S_1 and S_2 are vector subspaces of R^m , then

$$\dim(S_1 + S_2) = \dim(S_1) + \dim(S_2) - \dim(S_1 \cap S_2).$$

Example 2.19 Let S_1 and S_2 be subspaces of R^5 having bases $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ and $\{\mathbf{y}_1, \mathbf{y}_2\}$, respectively, where

$$\mathbf{x}_1 = (1, 0, 0, 1, 0)',$$

$$\mathbf{x}_2 = (0, 0, 1, 0, 1)',$$

$$\mathbf{x}_3 = (0, 1, 0, 0, 0)',$$

$$\mathbf{y}_1 = (1, 0, 0, 1, 1)',$$

$$\mathbf{y}_2 = (0, 1, 1, 0, 0)'.$$

We wish to find bases for $S_1 + S_2$ and $S_1 \cap S_2$. Now, clearly, $S_1 + S_2$ is spanned by the set $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{y}_1, \mathbf{y}_2\}$. Note that $\mathbf{y}_2 = \mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3 - \mathbf{y}_1$, and it can be easily verified that no constants $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ exist, except $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0$, which satisfy $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 + \alpha_3\mathbf{x}_3 + \alpha_4\mathbf{y}_1 = \mathbf{0}$. Thus, $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{y}_1\}$ is a basis for $S_1 + S_2$, and so $\dim(S_1 + S_2) = 4$. From Theorem 2.25, we know that

$\dim(S_1 \cap S_2) = 3 + 2 - 4 = 1$, and so any basis for $S_1 \cap S_2$ consists of one vector. The dependency between the \mathbf{x} 's and the \mathbf{y} 's will indicate an appropriate vector, so we seek solutions for $\alpha_1, \alpha_2, \alpha_3, \beta_1$, and β_2 , which satisfy

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \alpha_3 \mathbf{x}_3 = \beta_1 \mathbf{y}_1 + \beta_2 \mathbf{y}_2.$$

As a result, we find that a basis for $S_1 \cap S_2$ is given by the vector $\mathbf{y}_1 + \mathbf{y}_2 = (1, 1, 1, 1, 1)'$ because $\mathbf{x}_1 + \mathbf{x}_2 + \mathbf{x}_3 = \mathbf{y}_1 + \mathbf{y}_2$.

When S_1 and S_2 are such that $S_1 \cap S_2 = \{\mathbf{0}\}$, then the vector space obtained as the sum of S_1 and S_2 is sometimes referred to as the direct sum of S_1 and S_2 and written $S_1 \oplus S_2$. In this special case, each $\mathbf{x} \in S_1 \oplus S_2$ has a unique representation as $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$, where $\mathbf{x}_1 \in S_1$ and $\mathbf{x}_2 \in S_2$. A further special case is one in which S_1 and S_2 are orthogonal vector spaces; that is, for any $\mathbf{x}_1 \in S_1$ and $\mathbf{x}_2 \in S_2$, we have $\mathbf{x}_1' \mathbf{x}_2 = 0$. In this case, the unique representation $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$ for $\mathbf{x} \in S_1 \oplus S_2$ will have the vector \mathbf{x}_1 given by the orthogonal projection of \mathbf{x} onto S_1 , whereas \mathbf{x}_2 will be given by the orthogonal projection of \mathbf{x} onto S_2 . For instance, for any vector subspace S of R^m , $R^m = S \oplus S^\perp$, and for any $\mathbf{x} \in R^m$,

$$\mathbf{x} = P_S \mathbf{x} + P_{S^\perp} \mathbf{x}.$$

In general, if a vector space S is the sum of the r vector spaces S_1, \dots, S_r , and $S_i \cap S_j = \{\mathbf{0}\}$ for all $i \neq j$, then S is said to be the direct sum of S_1, \dots, S_r and is written as $S = S_1 \oplus \dots \oplus S_r$.

Example 2.20 Consider the vector spaces S_1, \dots, S_m , where S_i is spanned by $\{e_i\}$ and, as usual, e_i is the i th column of the $m \times m$ identity matrix. Consider a second sequence of vector spaces, T_1, \dots, T_m , where T_i is spanned by $\{e_i, e_{i+1}\}$ if $i \leq m-1$, whereas T_m is spanned by $\{e_1, e_m\}$. Then it follows that $R^m = S_1 + \dots + S_m$, as well as $R^m = T_1 + \dots + T_m$. However, although $R^m = S_1 \oplus \dots \oplus S_m$, it does not follow that $R^m = T_1 \oplus \dots \oplus T_m$, because it is not true that $T_i \cap T_j = \{\mathbf{0}\}$ for all $i \neq j$. Thus, any $\mathbf{x} = (x_1, \dots, x_m)'$ in R^m can be expressed uniquely as a sum comprised of a vector from each of the spaces S_1, \dots, S_m ; namely,

$$\mathbf{x} = x_1 e_1 + \dots + x_m e_m,$$

where $e_i \in S_i$. On the other hand, the decomposition corresponding to T_1, \dots, T_m is not unique. For instance, we can get the same sum above by choosing $e_1 \in T_1, e_2 \in T_2, \dots, e_m \in T_m$. However, we also have $\mathbf{x} = \mathbf{y}_1 + \dots + \mathbf{y}_m$ and $\mathbf{y}_i \in T_i$, where $\mathbf{y}_i = \frac{1}{2}(x_i e_i + x_{i+1} e_{i+1})$ when $i \leq m-1$ and $\mathbf{y}_m = \frac{1}{2}(x_1 e_1 + x_m e_m)$. In addition, the sum of the orthogonal projections of \mathbf{x} onto the spaces S_1, \dots, S_m yields \mathbf{x} , whereas the sum of the orthogonal projections of \mathbf{x} onto the spaces T_1, \dots, T_m yields $2\mathbf{x}$. Consider as a third sequence of vector spaces, U_1, \dots, U_m , where U_i has the basis $\{\gamma_i\}$ and $\gamma_i = e_1 + \dots + e_i$. Clearly, $U_i \cap U_j = \{\mathbf{0}\}$ if $i \neq j$, so $R^m = U_1 \oplus$

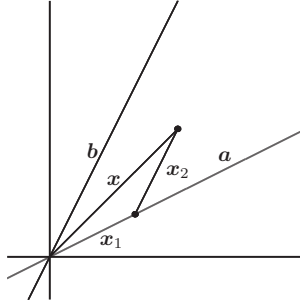


Figure 2.4 Projection of x onto a along b

$\cdots \oplus U_m$ and each $x \in R^m$ has a unique decomposition $x = x_1 + \cdots + x_m$ with $x_i \in U_i$. However, in this case, because the U_i 's are not orthogonal vector spaces, this decomposition of x is not given by the sum of the orthogonal projections of x onto the spaces U_1, \dots, U_m .

2.10 OBLIQUE PROJECTIONS

The projections that we have considered so far have been orthogonal projections. For any $x \in R^m$ and any subspace of R^m , S_1 , x can be uniquely expressed as $x = x_1 + x_2$, where $x_1 \in S_1$ and $x_2 \in S_2 = S_1^\perp$; x_1 is the orthogonal projection onto S_1 and x_2 is the orthogonal projection onto S_2 , and in this text these are the projections we are referring to when we simply use the term projection.

In this section, we generalize the idea of a projection to oblique projections, that is, to situations in which the subspaces S_1 and S_2 are not orthogonal to one another.

Definition 2.14 Let S_1 and S_2 be subspaces of R^m such that $S_1 \cap S_2 = \{\mathbf{0}\}$ and $S_1 \oplus S_2 = R^m$. Suppose

$$x = x_1 + x_2,$$

where $x \in R^m$, $x_1 \in S_1$, and $x_2 \in S_2$. Then x_1 is called the projection of x onto S_1 along S_2 , and x_2 is called the projection of x onto S_2 along S_1 .

An oblique projection in R^2 is depicted in Figure 2.4. Here S_1 is the line spanned by $a = (2, 1)'$ and S_2 is the line spanned by $b = (1, 2)'$. The projection of $x = (1, 1)'$ onto S_1 along S_2 is $x_1 = (2/3, 1/3)'$, while the projection onto S_2 along S_1 is $x_2 = (1/3, 2/3)'$.

When S_1 and S_2 are not orthogonal subspaces, we can find linear transformations of them such that these transformed subspaces are orthogonal. Suppose the columns of the $m \times m_1$ matrix X_1 form a basis for S_1 and the columns of the $m \times m_2$ matrix X_2 form a basis for S_2 . Then since $S_1 \cap S_2 = \{\mathbf{0}\}$ and $S_1 \oplus S_2 = R^m$, it follows that

$m_1 + m_2 = m$ and $X = (X_1, X_2)$ is nonsingular. For $i = 1, 2$, the subspace $T_i = \{\mathbf{y} : \mathbf{y} = X^{-1}\mathbf{x}, \mathbf{x} \in S_i\}$ has as a basis the columns of $Y_i = X^{-1}X_i$. Since $Y = (Y_1, Y_2) = (X^{-1}X_1, X^{-1}X_2) = X^{-1}X = I_m$, we have

$$Y_1 = \begin{bmatrix} I_{m_1} \\ (0) \end{bmatrix}, \quad Y_2 = \begin{bmatrix} (0) \\ I_{m_2} \end{bmatrix}, \quad (2.17)$$

and hence $Y_1'Y_2 = (0)$; that is, T_1 and T_2 are orthogonal subspaces and, in particular, T_2 is the orthogonal complement of T_1 . It follows from Theorem 2.17 that any $\mathbf{y} \in R^m$ can be expressed uniquely as $\mathbf{y} = \mathbf{y}_1 + \mathbf{y}_2$, where $\mathbf{y}_i \in T_i$. Equivalently, $\mathbf{x} = X\mathbf{y} = \mathbf{x}_1 + \mathbf{x}_2$, where $\mathbf{x}_i = X\mathbf{y}_i \in S_i$. This immediately leads to the following generalization of Theorem 2.17.

Theorem 2.26 Let S_1 and S_2 be subspaces of R^m such that $S_1 \cap S_2 = \{\mathbf{0}\}$ and $S_1 \oplus S_2 = R^m$. Then each $\mathbf{x} \in R^m$ can be expressed uniquely as

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2,$$

where $\mathbf{x}_i \in S_i$.

As with orthogonal projections, oblique projections can be computed through the use of a projection matrix. We will denote by $P_{S_1|S_2}$ the $m \times m$ matrix for which $P_{S_1|S_2}\mathbf{x}$ is the projection of \mathbf{x} onto S_1 along S_2 . Note that if \mathbf{y}_1 is the projection of $\mathbf{y} = X^{-1}\mathbf{x}$ onto T_1 along T_2 , that is, the orthogonal projection onto T_1 , then $X\mathbf{y}_1$ is the projection of \mathbf{x} onto S_1 along S_2 . We have seen in Section 2.7 that $\mathbf{y}_1 = P_{T_1}\mathbf{y} = Y_1(Y_1'Y_1)^{-1}Y_1'\mathbf{y}$, and so the required projection is $XY_1(Y_1'Y_1)^{-1}Y_1'X^{-1}\mathbf{x}$. Using (2.17), the projection matrix simplifies to the form given in the following theorem.

Theorem 2.27 Let S_1 and S_2 be subspaces of R^m such that $S_1 \cap S_2 = \{\mathbf{0}\}$ and $S_1 \oplus S_2 = R^m$. Suppose $X = (X_1, X_2)$, where the columns of the $m \times m_i$ matrix X_i form a basis for S_i . Then the projection matrix for the projection onto S_1 along S_2 is given by

$$P_{S_1|S_2} = X \begin{bmatrix} I_{m_1} & (0) \\ (0) & (0) \end{bmatrix} X^{-1}.$$

When S_1 and S_2 are orthogonal subspaces, $X_1'X_2 = (0)$ and

$$X^{-1} = \begin{bmatrix} (X_1'X_1)^{-1}X_1' \\ (X_2'X_2)^{-1}X_2' \end{bmatrix}.$$

In this case, $P_{S_1|S_2}$ as given in Theorem 2.27 reduces to $P_{S_1} = X_1(X_1'X_1)^{-1}X_1'$ as given in Section 2.7.

Example 2.21 For $i = 1, 2$, let S_i be the space spanned by the columns of X_i , where

$$X_1 = \begin{bmatrix} 1 & -1 \\ 1 & 1 \\ 2 & 1 \end{bmatrix}, \quad X_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \end{bmatrix}.$$

Since $X = (X_1, X_2)$ is nonsingular, $S_1 \cap S_2 = \{\mathbf{0}\}$ and $S_1 \oplus S_2 = R^3$. The projection matrix for the projection onto S_1 along S_2 is

$$\begin{aligned} P_{S_1|S_2} &= X \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} X^{-1} \\ &= \begin{bmatrix} 1 & -1 & 1 \\ 1 & 1 & 1 \\ 2 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 & -1 & 1 \\ -.5 & .5 & 0 \\ .5 & 1.5 & -1 \end{bmatrix} \\ &= \begin{bmatrix} .5 & -1.5 & 1 \\ -.5 & -.5 & 1 \\ -.5 & -1.5 & 2 \end{bmatrix}. \end{aligned}$$

The projection matrix for the projection onto S_2 along S_1 is then

$$P_{S_2|S_1} = I_3 - P_{S_1|S_2} = \begin{bmatrix} .5 & 1.5 & -1 \\ .5 & 1.5 & -1 \\ .5 & 1.5 & -1 \end{bmatrix}.$$

If $\mathbf{x} = (1, 2, 3)'$, then we find that $\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2$, where $\mathbf{x}_1 = P_{S_1|S_2} \mathbf{x} = (.5, 1.5, 2.5)' \in S_1$ and $\mathbf{x}_2 = P_{S_2|S_1} \mathbf{x} = (.5, .5, .5)' \in S_2$.

The projection matrix for an oblique projection, like that of an orthogonal projection, is idempotent since clearly $P_{S_1|S_2}^2 = P_{S_1|S_2}$. We saw in Section 2.7 that a matrix P is a projection matrix for an orthogonal projection if and only if P is idempotent and symmetric, implying then that $P'_{S_1|S_2} \neq P_{S_1|S_2}$ when S_1 and S_2 are not orthogonal subspaces. This is also evident from the form of $P_{S_1|S_2}$ given in Theorem 2.27 since if $P'_{S_1|S_2} = P_{S_1|S_2}$, then $X^{-1}P'_{S_1|S_2}X = X^{-1}P_{S_1|S_2}X$, or

$$(X'X)^{-1} \begin{bmatrix} I_{m_1} & (0) \\ (0) & (0) \end{bmatrix} (X'X) = \begin{bmatrix} I_{m_1} & (0) \\ (0) & (0) \end{bmatrix},$$

which holds only if $X'_1X_2 = (0)$.

We next consider another generalization of the projections discussed in Section 2.6. Instead of using the usual inner product, $\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}'\mathbf{y}$, we use the inner product $\langle \mathbf{x}, \mathbf{y} \rangle_A = \mathbf{x}'A\mathbf{y}$, where A is an $m \times m$ positive definite matrix. We refer to this as the A inner product.

Definition 2.15 Let S_1 be a subspace of R^m and suppose A is an $m \times m$ positive definite matrix. If

$$\mathbf{x} = \mathbf{x}_1 + \mathbf{x}_2,$$

where $\mathbf{x} \in R^m$, $\mathbf{x}_1 \in S_1$, and $\mathbf{x}_1' A \mathbf{x}_2 = 0$, then \mathbf{x}_1 is the orthogonal projection onto S_1 relative to the A inner product.

Let $U_1 = \{z : z = Bx, x \in S_1\}$, where B is any $m \times m$ matrix satisfying $B'B = A$. Then $z = z_1 + z_2$, where $z = Bx$, $z_1 = Bx_1 \in U_1$, and $z_1' z_2 = x_1' B' B x_2 = 0$, so that $z_2 = Bx_2 \in U_1^\perp$. Thus, the uniqueness of z_1 and z_2 guarantees the uniqueness of \mathbf{x}_1 and \mathbf{x}_2 in Definition 2.15. If the columns of X_1 form a basis for S_1 , then the columns of BX_1 form a basis for U_1 and the projection matrix for the orthogonal projection onto U_1 is $BX_1(X_1'AX_1)^{-1}X_1'B'$. From this, we find that $X_1(X_1'AX_1)^{-1}X_1'A$ is the projection matrix for the orthogonal projection onto S_1 relative to the A inner product.

It is not difficult to see that there is a direct connection between the decompositions given in Definitions 2.14 and 2.15. In Definition 2.14, if $X = (X_1, X_2)$ and the columns of X_i form a basis for S_i , then $x_1'X^{-1'}X^{-1}x_2 = 0$. Thus, \mathbf{x}_1 is the orthogonal projection onto S_1 and \mathbf{x}_2 is the orthogonal projection onto S_2 , both relative to the $(XX')^{-1}$ inner product. Similarly, in Definition 2.15, \mathbf{x}_1 is the projection onto S_1 along S_2 , where S_2 is the vector space which has the projection matrix $P_{S_2|S_1} = I_m - X_1(X_1'AX_1)^{-1}X_1'A$; S_2 is the orthogonal complement of S_1 , relative to the A inner product, since $X_1'AP_{S_2|S_1} = (0)$.

Suppose $\mathbf{x} \in R^m$ and we want to find a point in the subspace $S_1 \in R^m$ so that the distance between \mathbf{x} and that point, relative to the A inner product, is minimized. Using the decomposition in Definition 2.15 for \mathbf{x} , note that for any $\mathbf{y} \in S_1$, $\mathbf{x}_1 - \mathbf{y} \in S_1$, so that $(\mathbf{x} - \mathbf{x}_1)'A(\mathbf{x}_1 - \mathbf{y}) = \mathbf{x}_2'A(\mathbf{x}_1 - \mathbf{y}) = 0$. Thus,

$$\begin{aligned} (\mathbf{x} - \mathbf{y})'A(\mathbf{x} - \mathbf{y}) &= \{(\mathbf{x} - \mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{y})\}'A\{(\mathbf{x} - \mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{y})\} \\ &= (\mathbf{x} - \mathbf{x}_1)'A(\mathbf{x} - \mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{y})'A(\mathbf{x}_1 - \mathbf{y}) \\ &\quad + 2(\mathbf{x} - \mathbf{x}_1)'A(\mathbf{x}_1 - \mathbf{y}) \\ &= (\mathbf{x} - \mathbf{x}_1)'A(\mathbf{x} - \mathbf{x}_1) + (\mathbf{x}_1 - \mathbf{y})'A(\mathbf{x}_1 - \mathbf{y}) \\ &\geq (\mathbf{x} - \mathbf{x}_1)'A(\mathbf{x} - \mathbf{x}_1). \end{aligned}$$

That is, \mathbf{x}_1 is the point in S_1 that minimizes the distance, relative to the A inner product, from \mathbf{x} .

Example 2.22 In Example 2.17, we obtained the weighted least squares estimator of β in the multiple regression model

$$\mathbf{y} = X\beta + \epsilon,$$

where $\text{var}(\epsilon) = \sigma^2 \text{diag}(c_1^2, \dots, c_N^2)$, c_1^2, \dots, c_N^2 are known constants, and X has full column rank. We now consider a more general regression problem, sometimes

referred to as generalized least squares regression, in which $\text{var}(\epsilon) = \sigma^2 C$, where C is a known $N \times N$ positive definite matrix. Thus, the random errors not only may have different variances but also may be correlated, and weighted least squares regression is simply a special case of generalized least squares regression. When $C = I_N$ as it was in Example 2.11, we find $\hat{\beta}$ by minimizing the sum of squared errors

$$(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) = (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}),$$

where $\hat{\mathbf{y}} = X\hat{\beta}$. When $C \neq I_N$, we need to use the Mahalanobis distance function to measure the distance between \mathbf{y} and $X\hat{\beta}$. That is, to find the generalized least squares estimator $\hat{\beta}$, we need to minimize

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

Since $X\hat{\beta}$ is a point in the space spanned by the columns of X , and the closest point in this space to \mathbf{y} , relative to the C^{-1} inner product, is the orthogonal projection of \mathbf{y} onto this space, relative to the C^{-1} inner product, we must have

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

Thus, premultiplying by $(X'X)^{-1}X'$, we get

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

2.11 CONVEX SETS

A special type of subset of a vector space is known as a convex set. Such a set has the property that it contains any point on the line segment connecting any other two points in the set. Definition 2.16 follows.

Definition 2.16 A set $S \subseteq R^m$ is said to be a convex set if for any $\mathbf{x}_1 \in S$ and $\mathbf{x}_2 \in S$,

$$c\mathbf{x}_1 + (1 - c)\mathbf{x}_2 \in S,$$

where c is any scalar satisfying $0 < c < 1$.

The condition for a convex set is similar to the condition for a vector space; for S to be a vector space, we must have for any $\mathbf{x}_1 \in S$ and $\mathbf{x}_2 \in S$, $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 \in S$ for all α_1 and α_2 , whereas for S to be a convex set, this need only hold when α_1 and α_2 are nonnegative and $\alpha_1 + \alpha_2 = 1$. Thus, any vector space is a convex set. However, many familiar sets that are not vector spaces are, in fact, convex sets. For instance, intervals in R , rectangles in R^2 , and ellipsoidal regions in R^m are all examples of convex sets. The linear combination of \mathbf{x}_1 and \mathbf{x}_2 , $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2$, is called a convex

combination when $\alpha_1 + \alpha_2 = 1$ and $\alpha_i \geq 0$ for each i . More generally, $\alpha_1 \mathbf{x}_1 + \cdots + \alpha_r \mathbf{x}_r$ is called a convex combination of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ when $\alpha_1 + \cdots + \alpha_r = 1$ and $\alpha_i \geq 0$ for each i . Thus, by a simple induction argument, we see that a set S is convex if and only if it is closed under all convex combinations of vectors in S .

Theorem 2.28 indicates that the intersection of convex sets and the sum of convex sets are convex. The proof will be left as an exercise.

Theorem 2.28 Suppose that S_1 and S_2 are convex sets, where $S_i \subseteq R^m$ for each i . Then the set

- (a) $S_1 \cap S_2$ is convex,
- (b) $S_1 + S_2 = \{\mathbf{x}_1 + \mathbf{x}_2 : \mathbf{x}_1 \in S_1, \mathbf{x}_2 \in S_2\}$ is convex.

For any set S , the set $C(S)$ defined as the intersection of all convex sets containing S is called the convex hull of S . Consequently, because of a generalization of Theorem 2.28(a), $C(S)$ is the smallest convex set containing S .

A point \mathbf{a} is a limit or accumulation point of a set $S \subseteq R^m$ if for any $\delta > 0$, the set $S_\delta = \{\mathbf{x} : \mathbf{x} \in R^m, (\mathbf{x} - \mathbf{a})'(\mathbf{x} - \mathbf{a}) < \delta\}$ contains at least one point of S distinct from \mathbf{a} . A closed set is one that contains all of its limit points. If S is a set, then \bar{S} will denote its closure; that is, if S_0 is the set of all limit points of S , then $\bar{S} = S \cup S_0$. In Theorem 2.29, we see that the convexity of S guarantees the convexity of \bar{S} .

Theorem 2.29 If $S \subseteq R^m$ is a convex set, then its closure \bar{S} is also a convex set.

Proof. It is easily verified that the set $B_n = \{\mathbf{x} : \mathbf{x} \in R^m, \mathbf{x}'\mathbf{x} \leq n^{-1}\}$ is a convex set, where n is a positive integer. Consequently, it follows from Theorem 2.28(b) that $C_n = S + B_n$ is also convex. It also follows from a generalization of the result given in Theorem 2.28(a) that the set

$$A = \bigcap_{n=1}^{\infty} C_n$$

is convex. The result now follows by observing that $A = \bar{S}$. □

One of the most important results regarding convex sets is a theorem known as the separating hyperplane theorem. A hyperplane in R^m is a set of the form $T = \{\mathbf{x} : \mathbf{x} \in R^m, \mathbf{a}'\mathbf{x} = c\}$, where \mathbf{a} is an $m \times 1$ vector and c is a scalar. Thus, if $m = 2$, T represents a line in R^2 and if $m = 3$, T is a plane in R^3 . We will see that the separating hyperplane theorem states that two convex sets S_1 and S_2 are separated by a hyperplane if their intersection is empty; that is, a hyperplane exists that partitions R^m into two parts so that S_1 is contained in one part, whereas S_2 is contained in the other. Before proving this result, we will need to obtain some preliminary results. Our first result is a special case of the separating hyperplane theorem in which one of the sets contains the single point $\mathbf{0}$.

Theorem 2.30 Let S be a nonempty closed convex subset of R^m and suppose that $\mathbf{0} \notin S$. Then an $m \times 1$ vector \mathbf{a} exists, such that $\mathbf{a}'\mathbf{x} > 0$ for all $\mathbf{x} \in S$.

Proof. Let \mathbf{a} be a point in S which satisfies

$$\mathbf{a}'\mathbf{a} = \inf_{\mathbf{x} \in S} \mathbf{x}'\mathbf{x},$$

where \inf denotes the infimum or greatest lower bound. It is a consequence of the fact that S is closed and nonempty that such an $\mathbf{a} \in S$ exists. In addition, $\mathbf{a} \neq \mathbf{0}$ because $\mathbf{0} \notin S$. Now let c be an arbitrary scalar and \mathbf{x} any vector in S except for \mathbf{a} , and consider the vector $c\mathbf{x} + (1 - c)\mathbf{a}$. The squared length of this vector as a function of c is given by

$$\begin{aligned} f(c) &= \{c\mathbf{x} + (1 - c)\mathbf{a}\}'\{c\mathbf{x} + (1 - c)\mathbf{a}\} \\ &= \{c(\mathbf{x} - \mathbf{a}) + \mathbf{a}\}'\{c(\mathbf{x} - \mathbf{a}) + \mathbf{a}\} \\ &= c^2(\mathbf{x} - \mathbf{a})'(\mathbf{x} - \mathbf{a}) + 2c\mathbf{a}'(\mathbf{x} - \mathbf{a}) + \mathbf{a}'\mathbf{a}. \end{aligned}$$

Since the second derivative of this quadratic function $f(c)$ is positive, we find that it has a unique minimum at the point

$$c_* = -\frac{\mathbf{a}'(\mathbf{x} - \mathbf{a})}{(\mathbf{x} - \mathbf{a})'(\mathbf{x} - \mathbf{a})}.$$

Note that because S is convex, $\mathbf{x}_c = c\mathbf{x} + (1 - c)\mathbf{a} \in S$ when $0 \leq c \leq 1$, and so we must have $\mathbf{x}'_c \mathbf{x}_c = f(c) \geq f(0) = \mathbf{a}'\mathbf{a}$ for $0 \leq c \leq 1$ because of the way \mathbf{a} was defined. However, because of the quadratic structure of $f(c)$, this implies that $f(c) > f(0)$ for all $c > 0$. In other words, $c_* \leq 0$, which leads to

$$\mathbf{a}'(\mathbf{x} - \mathbf{a}) \geq 0$$

or

$$\mathbf{a}'\mathbf{x} \geq \mathbf{a}'\mathbf{a} > 0,$$

and this completes the proof. \square

A point \mathbf{x}_* is an interior point of S if a $\delta > 0$ exists, such that the set $S_\delta = \{\mathbf{x} : \mathbf{x} \in R^m, (\mathbf{x} - \mathbf{x}_*)'(\mathbf{x} - \mathbf{x}_*) < \delta\}$ is a subset of S . On the other hand, \mathbf{x}_* is a boundary point of S if for each $\delta > 0$, the set S_δ contains at least one point in S and at least one point not in S . Theorem 2.31 shows that the sets S and \bar{S} have the same interior points if S is convex.

Theorem 2.31 Suppose that S is a convex subset of R^m , whereas T is an open subset of R^m . If $T \subset \bar{S}$, then $T \subset S$.

Proof. Let \mathbf{x}_* be an arbitrary point in T , and define the sets

$$S_* = \{\mathbf{x} : \mathbf{x} = \mathbf{y} - \mathbf{x}_*, \mathbf{y} \in S\}, \quad T_* = \{\mathbf{x} : \mathbf{x} = \mathbf{y} - \mathbf{x}_*, \mathbf{y} \in T\}.$$

It follows from the conditions of Theorem 2.31 that S_* is convex, T_* is open, and $T_* \subset \bar{S}_*$. The proof will be complete if we can show that $\mathbf{0} \in S_*$ because this will imply that $\mathbf{x}_* \in S$. Since $\mathbf{0} \in T_*$ and T_* is an open set, we can find an $\epsilon > 0$ such that each of the vectors, $\epsilon \mathbf{e}_1, \dots, \epsilon \mathbf{e}_m, -\epsilon \mathbf{1}_m$ are in T_* . Since these vectors also must be in \bar{S}_* , we can find sequences, $\mathbf{x}_{i1}, \mathbf{x}_{i2}, \dots$, for $i = 1, 2, \dots, m+1$, such that each $\mathbf{x}_{ij} \in S_*$ and $\mathbf{x}_{ij} \rightarrow \epsilon \mathbf{e}_i$ for $i = 1, \dots, m$, and $\mathbf{x}_{ij} \rightarrow -\epsilon \mathbf{1}_m$ for $i = m+1$, as $j \rightarrow \infty$. Define the $m \times m$ matrix $X_j = (\mathbf{x}_{1j}, \dots, \mathbf{x}_{mj})$ so that $X_j \rightarrow \epsilon I_m$, as $j \rightarrow \infty$. It follows that an integer N_1 exists, such that X_j is nonsingular for all $j > N_1$. For $j > N_1$, define

$$\mathbf{y}_j = X_j^{-1} \mathbf{x}_{m+1,j} \quad (2.18)$$

so that

$$\mathbf{y}_j \rightarrow (\epsilon I_m)^{-1}(-\epsilon \mathbf{1}_m) = -\mathbf{1}_m.$$

Thus, some integer $N_2 \geq N_1$ exists, such that for all $j > N_2$, all of the components of \mathbf{y}_j are negative. But from (2.18), we have

$$\mathbf{x}_{m+1,j} - X_j \mathbf{y}_j = [X_j \quad \mathbf{x}_{m+1,j}] \begin{bmatrix} -\mathbf{y}_j \\ 1 \end{bmatrix} = \mathbf{0}.$$

This same equation holds if we replace the vector $(-\mathbf{y}'_j, 1)'$ by the unit vector $(\mathbf{y}'_j \mathbf{y}_j + 1)^{-1/2}(-\mathbf{y}'_j, 1)'$. Thus, $\mathbf{0}$ is a convex combination of the columns of $[X_j \quad \mathbf{x}_{m+1,j}]$, each of which is in S_* , so because S_* is convex, $\mathbf{0} \in S_*$. \square

Theorem 2.32 is sometimes called the supporting hyperplane theorem. It states that for any boundary point of a convex set S , a hyperplane passing through that point exists, such that none of the points of S are on one side of the hyperplane.

Theorem 2.32 Suppose that S is a convex subset of R^m and that \mathbf{x}_* either is not in S or is a boundary point of S if it is in S . Then an $m \times 1$ vector $\mathbf{b} \neq \mathbf{0}$ exists, such that $\mathbf{b}'\mathbf{x} \geq \mathbf{b}'\mathbf{x}_*$ for all $\mathbf{x} \in S$.

Proof. It follows from the previous theorem that \mathbf{x}_* also is not in \bar{S} or must be a boundary point of \bar{S} if it is in \bar{S} . Consequently, there exists a sequence of vectors, $\mathbf{x}_1, \mathbf{x}_2, \dots$ with each $\mathbf{x}_i \notin \bar{S}$, such that $\mathbf{x}_i \rightarrow \mathbf{x}_*$ as $i \rightarrow \infty$. Corresponding to each \mathbf{x}_i , define the set $S_i = \{\mathbf{y} : \mathbf{y} = \mathbf{x} - \mathbf{x}_i, \mathbf{x} \in S\}$, and note that $\mathbf{0} \notin \bar{S}_i$ because $\mathbf{x}_i \notin \bar{S}$. Thus, because \bar{S}_i is closed and convex by Theorem 2.29, it follows from Theorem 2.30 that an $m \times 1$ vector \mathbf{a}_i exists, such that $\mathbf{a}'_i \mathbf{y} > 0$ for all $\mathbf{y} \in \bar{S}_i$ or, equivalently, $\mathbf{a}'_i(\mathbf{x} - \mathbf{x}_i) > 0$ for all $\mathbf{x} \in \bar{S}$. Alternatively, we can write this as $\mathbf{b}'_i(\mathbf{x} - \mathbf{x}_i) > 0$, where $\mathbf{b}_i = (\mathbf{a}'_i \mathbf{a}_i)^{-1/2} \mathbf{a}_i$. Now because $\mathbf{b}'_i \mathbf{b}_i = 1$, the sequence $\mathbf{b}_1, \mathbf{b}_2, \dots$ is a

bounded sequence, and so it has a convergent subsequence; that is, positive integers $i_1 < i_2 < \cdots$ and some $m \times 1$ unit vector \mathbf{b} exist, such that $\mathbf{b}_{i_j} \rightarrow \mathbf{b}$ as $j \rightarrow \infty$. Consequently, $\mathbf{b}'_{i_j}(\mathbf{x} - \mathbf{x}_{i_j}) \rightarrow \mathbf{b}'(\mathbf{x} - \mathbf{x}_*)$ as $j \rightarrow \infty$, and we must have $\mathbf{b}'(\mathbf{x} - \mathbf{x}_*) \geq 0$ for all $\mathbf{x} \in S$ because $\mathbf{b}'_{i_j}(\mathbf{x} - \mathbf{x}_{i_j}) > 0$ for all $\mathbf{x} \in S$. This completes the proof. \square

We are now ready to prove the separating hyperplane theorem.

Theorem 2.33 Let S_1 and S_2 be convex subsets of R^m with $S_1 \cap S_2 = \emptyset$. Then an $m \times 1$ vector $\mathbf{b} \neq \mathbf{0}$ exists, such that $\mathbf{b}'\mathbf{x}_1 \geq \mathbf{b}'\mathbf{x}_2$ for all $\mathbf{x}_1 \in S_1$ and all $\mathbf{x}_2 \in S_2$.

Proof. Clearly, the set $S_{2*} = \{\mathbf{x} : -\mathbf{x} \in S_2\}$ is convex because S_2 is convex. Thus, from Theorem 2.28, we know that the set

$$S = S_1 + S_{2*} = \{\mathbf{x} : \mathbf{x} = \mathbf{x}_1 - \mathbf{x}_2, \mathbf{x}_1 \in S_1, \mathbf{x}_2 \in S_2\}$$

is also convex. In addition, $\mathbf{0} \notin S$ because $S_1 \cap S_2 = \emptyset$. Consequently, using Theorem 2.32, we find that an $m \times 1$ vector $\mathbf{b} \neq \mathbf{0}$ exists for which $\mathbf{b}'\mathbf{x} \geq 0$ for all $\mathbf{x} \in S$. However, this result implies that $\mathbf{b}'(\mathbf{x}_1 - \mathbf{x}_2) \geq 0$ for all $\mathbf{x}_1 \in S_1$ and all $\mathbf{x}_2 \in S_2$, as is required. \square

Suppose that $f(x)$ is a nonnegative function that is symmetric about $x = 0$ and has only one maximum, occurring at $x = 0$; in other words, $f(x) = f(-x)$ for all x and $f(x) \leq f(cx)$ if $0 \leq c \leq 1$. Clearly, the integral of $f(x)$ over an interval of fixed length will be maximized when the interval is centered at 0. This can be expressed as

$$\int_{-a}^a f(x + cy) \, dx \geq \int_{-a}^a f(x + y) \, dx,$$

for any $y, a > 0$, and $0 \leq c \leq 1$. This result has some important applications regarding probabilities of random variables. The following result, which is a generalization to a function $f(\mathbf{x})$ of the $m \times 1$ vector \mathbf{x} replaces the interval in R^1 by a symmetric convex set in R^m . This generalization is due to Anderson (1955). For some simple applications of the result to probabilities of random vectors, see Problem 2.68. For additional extensions and applications of this result, see Anderson (1996) and Perlman (1990).

Theorem 2.34 Let S be a convex subset of R^m , symmetric about $\mathbf{0}$, so that if $\mathbf{x} \in S$, $-\mathbf{x} \in S$ also. Let $f(\mathbf{x}) \geq 0$ be a function for which $f(\mathbf{x}) = f(-\mathbf{x})$, $S_\alpha = \{\mathbf{x} : f(\mathbf{x}) \geq \alpha\}$ is convex for any positive α , and $\int_S f(\mathbf{x}) \, d\mathbf{x} < \infty$. Then

$$\int_S f(\mathbf{x} + c\mathbf{y}) \, d\mathbf{x} \geq \int_S f(\mathbf{x} + \mathbf{y}) \, d\mathbf{x},$$

for $0 \leq c \leq 1$ and $\mathbf{y} \in R^m$.

A more comprehensive discussion of convex sets can be found in Berkovitz (2002), Kelly and Weiss (1979), Lay (1982), and Rockafellar (1970), whereas some applications of the separating hyperplane theorem to statistical decision theory can be found in Ferguson (1967).

PROBLEMS

2.1 Determine whether each of the following sets of vectors is a vector space:

- (a) $\{(a, b, a + b, 1)' : -\infty < a < \infty, -\infty < b < \infty\}$.
- (b) $\{(a, b, c, a + b - 2c)' : -\infty < a < \infty, -\infty < b < \infty, -\infty < c < \infty\}$.
- (c) $\{(a, b, c, 1 - a - b - c)' : -\infty < a < \infty, -\infty < b < \infty, -\infty < c < \infty\}$.

2.2 Consider the vector space

$$S = \{(a, a + b, a + b, -b)' : -\infty < a < \infty, -\infty < b < \infty\}.$$

Determine which of the following sets of vectors are spanning sets of S :

- (a) $\{(1, 0, 0, 1)', (1, 2, 2, -1)'\}$.
- (b) $\{(1, 1, 0, 0)', (0, 0, 1, -1)'\}$.
- (c) $\{(2, 1, 1, 1)', (3, 1, 1, 2)', (3, 2, 2, 1)'\}$.
- (d) $\{(1, 0, 0, 0)', (0, 1, 1, 0)', (0, 0, 0, 1)'\}$.

2.3 Is the vector $\mathbf{x} = (1, 1, 1, 1)'$ in the vector space S given in Problem 2.2? Is the vector $\mathbf{y} = (4, 1, 1, 3)'$ in S ?

2.4 Let $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ be a set of vectors in a vector space S , and let W be the vector subspace consisting of all possible linear combinations of these vectors. Prove that W is the smallest subspace of S that contains the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$; that is, show that if V is another vector subspace containing $\mathbf{x}_1, \dots, \mathbf{x}_r$, then W is a subspace of V .

2.5 Suppose that \mathbf{x} is a random vector having a distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix Ω given by

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \Omega = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}.$$

Let $\mathbf{x}_1 = (2, 2)'$ and $\mathbf{x}_2 = (2, 0)'$ be two observations from this distribution. Use the Mahalanobis distance function to determine which of these two observations is closer to the mean.

2.6 Use the Cauchy–Schwarz inequality to prove the triangle inequality; that is, for $m \times 1$ vectors \mathbf{x} and \mathbf{y} , show that

$$\{(\mathbf{x} + \mathbf{y})'(\mathbf{x} + \mathbf{y})\}^{1/2} \leq (\mathbf{x}'\mathbf{x})^{1/2} + (\mathbf{y}'\mathbf{y})^{1/2}.$$

2.7 Using the law of cosines, find

(a) the angle between $\mathbf{x} = (1, 2, 1, 2)'$ and $\mathbf{y} = (3, 0, 1, 1)'$,

(b) $\mathbf{x}'\mathbf{y}$ if $\mathbf{x}'\mathbf{x} = 3$, $\mathbf{y}'\mathbf{y} = 2$, and the angle between \mathbf{x} and \mathbf{y} is $\pi/6$.

2.8 Show that the functions $\|\mathbf{x}\|_p$ and $\|\mathbf{x}\|_\infty$ defined in Section 2.2 are, in fact, vector norms.

2.9 If a vector norm is derived from an inner product, that is, $\|\mathbf{x}\| = \langle \mathbf{x}, \mathbf{x} \rangle^{1/2}$, show that the identity

$$\|\mathbf{x} + \mathbf{y}\|^2 + \|\mathbf{x} - \mathbf{y}\|^2 = 2\|\mathbf{x}\|^2 + 2\|\mathbf{y}\|^2,$$

known as the parallelogram identity, holds for all \mathbf{x} and \mathbf{y} .

2.10 Prove Theorem 2.3.

2.11 Which of the following sets of vectors are linearly dependent?

(a) $\{(1, -1, 2)', (3, 1, 1)'\}$.

(b) $\{(4, -1, 2)', (3, 2, 3)', (2, 5, 4)'\}$.

(c) $\{(1, 2, 3)', (2, 3, 1)', (-1, 1, 1)'\}$.

(d) $\{(1, -1, -1)', (2, 4, 3)', (3, 3, 5)', (7, 0, -1)'\}$.

2.12 Show that the set of vectors $\{(1, 2, 2, 2)', (1, 2, 1, 2)', (1, 1, 1, 1)'\}$ is a linearly independent set.

2.13 Consider the set of vectors

$$\{(2, 1, 4, 3)', (3, 0, 5, 2)', (0, 3, 2, 5)', (4, 2, 8, 6)'\}.$$

(a) Show that this set of vectors is linearly dependent.

(b) From this set of four vectors find a subset of two vectors that is a linearly independent set.

2.14 Prove that the set of $m \times 1$ vectors $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$ is a linearly independent set if and only if $\{\mathbf{x}_1, \mathbf{x}_1 + \mathbf{x}_2, \dots, \sum_{i=1}^n \mathbf{x}_i\}$ is a linearly independent set.

2.15 Which of the following sets of vectors are bases for R^4 ?

(a) $\{(0, 1, 0, 1)', (1, 1, 0, 0)', (0, 0, 1, 1)'\}$.

(b) $\{(2, 2, 2, 1)', (2, 1, 1, 1)', (3, 2, 1, 1)', (1, 1, 1, 1)'\}$.

(c) $\{(2, 0, 1, 1)', (3, 1, 2, 2)', (2, 1, 1, 2)', (2, 1, 2, 1)'\}$.

2.16 Let $\mathbf{x}_1 = (2, -3, 2)'$ and $\mathbf{x}_2 = (4, 1, 1)'$. Find a vector \mathbf{x}_3 so that the set $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$ forms a basis for R^3 .

- 2.17** Suppose that the vector space S is spanned by the vectors $\mathbf{x}_1 = (1, -2, 1)'$, $\mathbf{x}_2 = (2, 1, 1)'$, and $\mathbf{x}_3 = (8, -1, 5)'$.
- (a) Show that the dimension of S is two and find a basis, $\{\mathbf{z}_1, \mathbf{z}_2\}$, for S .
 - (b) Show that the vector $\mathbf{x} = (1, 3, 0)'$ is in S by finding the scalars α_1 and α_2 , such that $\mathbf{x} = \alpha_1 \mathbf{z}_1 + \alpha_2 \mathbf{z}_2$.
 - (c) For the \mathbf{x} given in part (b), find two different sets of scalars $\alpha_1, \alpha_2, \alpha_3$, such that $\mathbf{x} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \alpha_3 \mathbf{x}_3$.
- 2.18** Show that if $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ is a basis for a vector space S , then every set containing more than r vectors from S must be linearly dependent. Note that this result proves that the number of vectors in a basis for S is uniquely defined.
- 2.19** Prove the results of Theorem 2.9.
- 2.20** Prove that if a set of orthogonal vectors does not contain the null vector, then it is a linearly independent set.
- 2.21** Suppose that \mathbf{x}_1 and \mathbf{x}_2 are linearly independent and define $\mathbf{y}_1 = a\mathbf{x}_1 + b\mathbf{x}_2$ and $\mathbf{y}_2 = c\mathbf{x}_1 + d\mathbf{x}_2$. Show that \mathbf{y}_1 and \mathbf{y}_2 are linearly independent if and only if $ad \neq bc$.
- 2.22** Find a basis for the vector space given in Problem 2.2. What is the dimension of this vector space? Find a second different basis for this same vector space, where none of the vectors in this second basis is a scalar multiple of a vector in the first basis.
- 2.23** Show that the set of vectors $\{\gamma_1, \dots, \gamma_m\}$, given in Example 2.4, is a basis for R^m .
- 2.24** Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that
- (a) $R(AB) \subseteq R(A)$.
 - (b) $R(AB) = R(A)$ if $\text{rank}(AB) = \text{rank}(A)$.
- 2.25** Suppose A and B are $m \times n$ matrices. Show that an $n \times n$ matrix C exists, which satisfies $AC = B$ if and only if $R(B) \subseteq R(A)$.
- 2.26** Let A be an $m \times n$ matrix and B be an $m \times p$ matrix. Prove that

$$\text{rank}([A \ B]) \leq \text{rank}(A) + \text{rank}(B).$$

- 2.27** Prove the results of Theorem 2.14.
- 2.28** Suppose \mathbf{x} is an $m \times 1$ nonnull vector and \mathbf{y} is an $n \times 1$ nonnull vector. What is the rank of the matrix $\mathbf{x}\mathbf{y}'$?
- 2.29** Let A , B , and C be $p \times n$, $m \times q$, and $m \times n$ matrices, respectively. Prove that

$$\text{rank} \left(\begin{bmatrix} C & B \\ A & (0) \end{bmatrix} \right) = \text{rank}(A) + \text{rank}(B)$$

if an $m \times p$ matrix F and a $q \times n$ matrix G exist, such that $C = FA + BG$.

- 2.30** Let A be an $m \times n$ matrix and B be an $n \times p$ matrix with $\text{rank}(B) = n$. Show that $\text{rank}(A) = \text{rank}(AB)$.

- 2.31** Show by example that if A and B are $m \times m$ matrices, then $\text{rank}(AB) = \text{rank}(BA)$ does not necessarily hold.
- 2.32** Refer to Example 2.7 and Example 2.10. Find the matrix A satisfying $Z_1 = X_1 A$, where $Z_1 = (z_1, z_2, z_3)$ and $X_1 = (x_1, x_2, x_3)$. Show that $AA' = (X_1' X_1)^{-1}$.
- 2.33** Let S be the vector space spanned by the vectors $x_1 = (1, 2, 1, 2)'$, $x_2 = (2, 3, 1, 2)'$, $x_3 = (3, 4, -1, 0)'$, and $x_4 = (3, 4, 0, 1)'$.
- Find a basis for S .
 - Use the Gram–Schmidt procedure on the basis found in (a) to determine an orthonormal basis for S .
 - Find the orthogonal projection of $x = (1, 0, 0, 1)'$ onto S .
 - Find the component of x orthogonal to S .
- 2.34** Using (2.10), determine the projection matrix for the vector space S given in Problem 2.33. Use the projection matrix to compute the orthogonal projection of $x = (1, 0, 0, 1)'$ onto S .
- 2.35** Let S be the vector space spanned by the vectors $x_1 = (1, 2, 3)'$ and $x_2 = (1, 1, -1)'$. Find the point in S that is closest to the point $x = (1, 1, 1)'$.
- 2.36** Let $\{z_1, \dots, z_r\}$ be an orthonormal basis for a vector space S . Show that if $x \in S$, then

$$x'x = (x'z_1)^2 + \dots + (x'z_r)^2.$$

- 2.37** Suppose S is a vector subspace of R^4 having the projection matrix

$$P_S = \frac{1}{10} \begin{bmatrix} 6 & -2 & -2 & -4 \\ -2 & 9 & -1 & -2 \\ -2 & -1 & 9 & -2 \\ -4 & -2 & -2 & 6 \end{bmatrix}.$$

- What is the dimension of S ?
 - Find a basis for S .
- 2.38** Let P_1 and P_2 be the projection matrices for the orthogonal projections onto $S_1 \in R^m$ and $S_2 \in R^m$, respectively. Show that
- $P_1 + P_2$ is a projection matrix if and only if $P_1 P_2 = P_2 P_1 = (0)$,
 - $P_1 - P_2$ is a projection matrix if and only if $P_1 P_2 = P_2 P_1 = P_2$.
- 2.39** Consider the vector space $S = \{u : u = Ax, x \in R^4\}$, where A is the 4×4 matrix given by

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

- (a) Determine the dimension of S and find a basis.
 (b) Determine the dimension of the null space $N(A)$ and find a basis for it.
 (c) Is the vector $(3, 5, 2, 4)'$ in S ?
 (d) Is the vector $(1, 1, 1, 1)'$ in $N(A)$?
- 2.40** Let $\mathbf{x} \in R^n$ and suppose that $\mathbf{u}(\mathbf{x})$ defines a linear transformation of R^n into R^m . Using the standard basis $\{\mathbf{e}_1, \dots, \mathbf{e}_n\}$ for R^n and the $m \times 1$ vectors $\mathbf{u}(\mathbf{e}_1), \dots, \mathbf{u}(\mathbf{e}_n)$, prove that an $m \times n$ matrix A exists, for which

$$\mathbf{u}(\mathbf{x}) = A\mathbf{x},$$

for every $\mathbf{x} \in R^n$.

- 2.41** Let T be a vector subspace of R^n , and suppose that S is the subspace of R^m given by

$$S = \{\mathbf{u}(\mathbf{x}) : \mathbf{x} \in T\},$$

where the transformation defined by \mathbf{u} is linear. Show that an $m \times n$ matrix A exists, which satisfies

$$\mathbf{u}(\mathbf{x}) = A\mathbf{x},$$

for every $\mathbf{x} \in T$.

- 2.42** Let T be the vector space spanned by the two vectors $\mathbf{x}_1 = (1, 1, 0)'$ and $\mathbf{x}_2 = (0, 1, 1)'$. Let S be the vector space defined as $S = \{\mathbf{u}(\mathbf{x}) : \mathbf{x} \in T\}$, where the function \mathbf{u} defines a linear transformation satisfying $\mathbf{u}(\mathbf{x}_1) = (2, 3, 1)'$ and $\mathbf{u}(\mathbf{x}_2) = (2, 5, 3)'$. Find a matrix A , such that $\mathbf{u}(\mathbf{x}) = A\mathbf{x}$, for all $\mathbf{x} \in T$.
- 2.43** Consider the linear transformation defined by

$$\mathbf{u}(\mathbf{x}) = \begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_m - \bar{x} \end{bmatrix},$$

for all $\mathbf{x} \in R^m$, where $\bar{x} = (1/m) \sum x_i$. Find the matrix A for which $\mathbf{u}(\mathbf{x}) = A\mathbf{x}$, and then determine the dimension of the range and null spaces of A .

- 2.44** In an introductory statistics course, students must take three 100-point exams followed by a 150-point final exam. We will identify the scores on these exams with the variables x_1, x_2, x_3 , and y . We want to be able to estimate the value of y once x_1, x_2 , and x_3 are known. A class of 32 students produced the following set of exam scores.

Student	x_1	x_2	x_3	y	Student	x_1	x_2	x_3	y
1	87	89	92	111	17	72	76	96	116
2	72	85	77	99	18	73	70	52	78
3	67	79	54	82	19	73	61	86	101
4	79	71	68	136	20	73	83	76	82
5	60	67	53	73	21	97	99	97	141
6	83	84	92	107	22	84	92	86	112
7	82	88	76	106	23	82	68	73	62
8	87	68	91	128	24	61	59	77	56
9	88	66	65	95	25	78	73	81	137
10	62	68	63	108	26	84	73	68	118
11	100	100	100	142	27	57	47	71	108
12	87	82	80	89	28	87	95	84	121
13	72	94	76	109	29	62	29	66	71
14	86	92	98	140	30	77	82	81	123
15	85	82	62	117	31	52	66	71	102
16	62	50	71	102	32	95	99	96	130

- (a) Find the least squares estimator for $\beta = (\beta_0, \beta_1, \beta_2, \beta_3)'$ in the multiple regression model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \epsilon.$$

- (b) Find the least squares estimator for $\beta_1 = (\beta_0, \beta_1, \beta_2)'$ in the model

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \epsilon.$$

- (c) Compute the reduction in the sum of squared errors attributable to the inclusion of the variable x_3 in the model given in (a).

2.45 Suppose that we have independent samples of a response y corresponding to k different treatments with a sample size of n_i responses from the i th treatment. If the j th observation from the i th treatment is denoted y_{ij} , then the model

$$y_{ij} = \mu_i + \epsilon_{ij}$$

is known as the one-way classification model. Here μ_i represents the expected value of a response from treatment i , whereas the ϵ_{ij} 's are independent and identically distributed as $N(0, \sigma^2)$.

- (a) If we let $\beta = (\mu_1, \dots, \mu_k)'$, write this model in matrix form by defining \mathbf{y} , \mathbf{X} , and $\boldsymbol{\epsilon}$ so that $\mathbf{y} = \mathbf{X}\beta + \boldsymbol{\epsilon}$.
- (b) Find the least squares estimator of β , and show that the sum of squared errors for the corresponding fitted model is given by

$$\text{SSE}_1 = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2,$$

where

$$\bar{y}_i = \sum_{j=1}^{n_i} y_{ij} / n_i.$$

(c) If $\mu_1 = \cdots = \mu_k = \mu$, then the reduced model

$$y_{ij} = \mu + \epsilon_{ij}$$

holds for all i and j . Find the least squares estimator of μ and the sum of squared errors SSE_2 for the fitted reduced model. Show that $\text{SSE}_2 - \text{SSE}_1$, referred to as the sum of squares for treatment and denoted SST, can be expressed as

$$\text{SST} = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2,$$

where

$$\bar{y} = \sum_{i=1}^k n_i \bar{y}_i / n, \quad n = \sum_{i=1}^k n_i.$$

(d) Show that the F statistic given in (2.11) takes the form

$$F = \frac{\text{SST}/(k-1)}{\text{SSE}_1/(n-k)}.$$

2.46 Suppose that we have the model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$ and wish to find the estimator $\hat{\boldsymbol{\beta}}$, which minimizes

$$(\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}),$$

subject to the restriction that $\hat{\boldsymbol{\beta}}$ satisfies $A\hat{\boldsymbol{\beta}} = \mathbf{0}$, where X has full column rank and A has full row rank.

(a) Show that $S = \{\mathbf{y} : \mathbf{y} = X\hat{\boldsymbol{\beta}}, A\hat{\boldsymbol{\beta}} = \mathbf{0}\}$ is a vector space.

(b) Let C be any matrix whose columns form a basis for the null space of A ; that is, C satisfies the identity, $C(C'C)^{-1}C' = I - A'(AA')^{-1}A$. Using the geometrical properties of least squares estimators, show that the restricted least squares estimator $\hat{\boldsymbol{\beta}}$ is given by

$$\hat{\boldsymbol{\beta}} = C(C'X'XC)^{-1}C'X'\mathbf{y}.$$

2.47 Let S_1 and S_2 be vector subspaces of R^m . Show that $S_1 + S_2$ also must be a vector subspace of R^m .

2.48 Let S_1 be the vector space spanned by $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$, where $\mathbf{x}_1 = (1, 1, 1, 1)'$, $\mathbf{x}_2 = (1, 2, 2, 2)'$, and $\mathbf{x}_3 = (1, 0, -2, -2)'$. Let S_2 be the vector space spanned by $\{\mathbf{y}_1, \mathbf{y}_2, \mathbf{y}_3\}$, where $\mathbf{y}_1 = (1, 3, 5, 5)'$, $\mathbf{y}_2 = (1, 2, 3, 6)'$, and $\mathbf{y}_3 = (0, 1, 4, 7)'$. Find bases for $S_1 + S_2$ and $S_1 \cap S_2$.

- 2.49** Let S_1 and S_2 be vector subspaces of R^m . Show that $S_1 + S_2$ is the vector space of smallest dimension containing $S_1 \cup S_2$. In other words, show that if T is a vector space for which $S_1 \cup S_2 \subseteq T$, then $S_1 + S_2 \subseteq T$.
- 2.50** Prove Theorem 2.25.
- 2.51** Let S_1 and S_2 be vector subspaces of R^m . Suppose that $\{x_1, \dots, x_r\}$ spans S_1 and $\{y_1, \dots, y_h\}$ spans S_2 . Show that $\{x_1, \dots, x_r, y_1, \dots, y_h\}$ spans the vector space $S_1 + S_2$.
- 2.52** Let S_1 be the vector space spanned by the vectors

$$x_1 = \begin{bmatrix} 3 \\ 1 \\ 3 \\ 1 \end{bmatrix}, \quad x_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad x_3 = \begin{bmatrix} 2 \\ 1 \\ 2 \\ 1 \end{bmatrix},$$

whereas the vector space S_2 is spanned by the vectors

$$y_1 = \begin{bmatrix} 3 \\ 0 \\ 5 \\ -1 \end{bmatrix}, \quad y_2 = \begin{bmatrix} 1 \\ 2 \\ 3 \\ 1 \end{bmatrix}, \quad y_3 = \begin{bmatrix} 1 \\ -4 \\ -1 \\ -3 \end{bmatrix}.$$

Find the following:

- (a) Bases for S_1 and S_2 .
 - (b) The dimension of $S_1 + S_2$.
 - (c) A basis for $S_1 + S_2$.
 - (d) The dimension of $S_1 \cap S_2$.
 - (e) A basis for $S_1 \cap S_2$.
- 2.53** Let S_1 and S_2 be vector subspaces of R^m with $\dim(S_1) = r_1$ and $\dim(S_2) = r_2$.
- (a) Find expressions in terms of m , r_1 , and r_2 for the smallest and largest possible values of $\dim(S_1 + S_2)$.
 - (b) Find the smallest and largest possible values of $\dim(S_1 \cap S_2)$.
- 2.54** Let S_1 and S_2 be vector subspaces of R^m .
- (a) Show that $(S_1 + S_2)^\perp = S_1^\perp \cap S_2^\perp$.
 - (b) Show that $(S_1 \cap S_2)^\perp = S_1^\perp + S_2^\perp$.
- 2.55** Let T be the vector space spanned by the vectors $\{(1, 1, 1)', (2, 1, 2)'\}$. Find a vector space S_1 , such that $R^3 = T \oplus S_1$. Find another vector space S_2 , such that $R^3 = T \oplus S_2$ and $S_1 \cap S_2 = \{0\}$.
- 2.56** Let S_1 be the vector space spanned by $\{(1, 1, -2, 0)', (2, 0, 1, -3)'\}$, whereas S_2 is spanned by $\{(1, 1, 1, -3)', (1, 1, 1, 1)'\}$. Show that $R^4 = S_1 + S_2$. Is this a direct sum? That is, can we write $R^4 = S_1 \oplus S_2$? Are S_1 and S_2 orthogonal vector spaces?

2.57 Let S_1 and S_2 be vector subspaces of R^m , and let $T = S_1 + S_2$. Show that this sum is a direct sum, that is, $T = S_1 \oplus S_2$ if and only if

$$\dim(T) = \dim(S_1) + \dim(S_2).$$

2.58 For $i = 1, 2$, let S_i be the space spanned by the columns of X_i , where

$$X_1 = \begin{bmatrix} 1 & 0 \\ -1 & 1 \\ -1 & -2 \end{bmatrix}, \quad X_2 = \begin{bmatrix} 1 \\ -1 \\ 0 \end{bmatrix}.$$

- (a) Find the projection matrix for the projection onto S_1 along S_2 .
- (b) Find the projection of $\mathbf{x} = (2, 2, 1)'$ onto S_1 along S_2 and the projection of \mathbf{x} onto S_2 along S_1 .

2.59 Consider the projection matrix

$$P_{S_1|S_2} = \begin{bmatrix} 1 & -1 & 1 \\ 0 & -2 & 3 \\ 0 & -2 & 3 \end{bmatrix}.$$

- (a) Find a basis for the vector space S_1 .
 - (b) Find a basis for the vector space S_2 .
- 2.60** For $i = 1, 2$, let S_i and T_i be subspaces of R^m such that $S_i \oplus T_i = R^m$ and $S_i \cap T_i = \{\mathbf{0}\}$. Show that $S_1 = S_2$ if and only if $P_{S_1|T_1}P_{S_2|T_2} = P_{S_2|T_2}$ and $P_{S_2|T_2}P_{S_1|T_1} = P_{S_1|T_1}$.
- 2.61** For $i = 1, 2$, suppose S_i and T_i are subspaces of R^m satisfying $S_1 \oplus S_2 = T_1 \oplus T_2 = R^m$ and $S_1 \cap S_2 = T_1 \cap T_2 = \{\mathbf{0}\}$.
- (a) Show that $P_{S_1|S_2} + P_{T_1|T_2}$ is a projection matrix if and only if $P_{S_1|S_2}P_{T_1|T_2} = P_{T_1|T_2}P_{S_1|S_2} = (\mathbf{0})$.
 - (b) When the condition in (a) holds, show that $P_{S_1|S_2} + P_{T_1|T_2}$ is the projection matrix for the projection onto $S_1 \oplus T_1$ along $S_2 \cap T_2$.
- 2.62** Determine which of the following sets are convex and which are not convex:
- (a) $S_1 = \{(x_1, x_2)' : x_1^2 + x_2^2 \leq 1\}$.
 - (b) $S_2 = \{(x_1, x_2)' : x_1^2 + x_2^2 = 1\}$.
 - (c) $S_3 = \{(x_1, x_2)' : 0 \leq x_1 \leq x_2 \leq 1\}$.
- 2.63** Prove Theorem 2.28.
- 2.64** Show that if S_1 and S_2 are convex subsets of R^m , then $S_1 \cup S_2$ need not be convex.
- 2.65** Show that for any positive scalar n , the set $B_n = \{\mathbf{x} : \mathbf{x} \in R^m, \mathbf{x}'\mathbf{x} \leq n^{-1}\}$ is convex.
- 2.66** For any set $S \subseteq R^m$, show that its convex hull $C(S)$ consists of all convex combinations of the vectors in S .

- 2.67** Suppose that S is a nonempty subset of R^m . Show that every vector in the convex hull of S can be expressed as a convex combination of $m + 1$ or fewer vectors in S .
- 2.68** Let \mathbf{x} be an $m \times 1$ random vector with density function $f(\mathbf{x})$ such that $f(\mathbf{x}) = f(-\mathbf{x})$ and the set $\{\mathbf{x} : f(\mathbf{x}) \geq \alpha\}$ is convex for all positive α . Suppose that S is a convex subset of R^m , symmetric about $\mathbf{0}$.
- (a) Show that $P(\mathbf{x} + c\mathbf{y} \in S) \geq P(\mathbf{x} + \mathbf{y} \in S)$ for any constant vector \mathbf{y} and $0 \leq c \leq 1$.
 - (b) Show that the inequality in (a) also holds if \mathbf{y} is an $m \times 1$ random vector distributed independently of \mathbf{x} .
 - (c) Show that if $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, its density function satisfies the conditions of this exercise.
 - (d) Show that if \mathbf{x} and \mathbf{y} are independently distributed with $\mathbf{x} \sim N_m(\mathbf{0}, \Omega_1)$ and $\mathbf{y} \sim N_m(\mathbf{0}, \Omega_2)$, such that $\Omega_1 - \Omega_2$ is nonnegative definite, then $P(\mathbf{x} \in S) \leq P(\mathbf{y} \in S)$.

3

EIGENVALUES AND EIGENVECTORS

3.1 INTRODUCTION

Eigenvalues and eigenvectors are special implicitly defined functions of the elements of a square matrix. In many applications involving the analysis of a square matrix, the key information from the analysis is provided by some or all of these eigenvalues and eigenvectors, which is a consequence of some of the properties of eigenvalues and eigenvectors that we will develop in this chapter. However, before we get to these properties, we must first understand how eigenvalues and eigenvectors are defined and how they are calculated.

3.2 EIGENVALUES, EIGENVECTORS, AND EIGENSPACES

If A is an $m \times m$ matrix, then any scalar λ satisfying the equation

$$Ax = \lambda x, \tag{3.1}$$

for some $m \times 1$ vector $x \neq \mathbf{0}$, is called an eigenvalue of A . The vector x is called an eigenvector of A corresponding to the eigenvalue λ , and (3.1) is called the eigenvalue-eigenvector equation of A . Eigenvalues and eigenvectors are also

sometimes referred to as latent roots and vectors or characteristic roots and vectors. Equation (3.1) can be equivalently expressed as

$$(A - \lambda I_m)x = 0. \quad (3.2)$$

Note that if $|A - \lambda I_m| \neq 0$, then $(A - \lambda I_m)^{-1}$ would exist, and so premultiplication of (3.2) by this inverse would lead to a contradiction of the already stated assumption that $x \neq 0$. Thus, any eigenvalue λ must satisfy the determinantal equation

$$|A - \lambda I_m| = 0,$$

which is known as the characteristic equation of A . Applying the definition of a determinant, we readily observe that the characteristic equation is an m th degree polynomial in λ ; that is, scalars $\alpha_1, \dots, \alpha_{m-1}$ exist, such that the characteristic equation above can be expressed equivalently as

$$(-\lambda)^m + \alpha_{m-1}(-\lambda)^{m-1} + \dots + \alpha_1(-\lambda) + \alpha_0 = 0.$$

Since an m th degree polynomial has m roots, it follows that an $m \times m$ matrix has m eigenvalues; that is, m scalars $\lambda_1, \dots, \lambda_m$ exist that satisfy the characteristic equation. When all of the eigenvalues of A are real, we will sometimes find it notationally convenient to identify the i th largest eigenvalue of the matrix A as $\lambda_i(A)$. In other words, in this case, the ordered eigenvalues of A may be written as $\lambda_1(A) \geq \dots \geq \lambda_m(A)$.

The characteristic equation can be used to obtain the eigenvalues of the matrix A . These can be then used in the eigenvalue-eigenvector equation to obtain corresponding eigenvectors.

Example 3.1 We will find the eigenvalues and eigenvectors of the 3×3 matrix A given by

$$A = \begin{bmatrix} 5 & -3 & 3 \\ 4 & -2 & 3 \\ 4 & -4 & 5 \end{bmatrix}.$$

The characteristic equation of A is

$$\begin{aligned} |A - \lambda I_3| &= \begin{vmatrix} 5 - \lambda & -3 & 3 \\ 4 & -2 - \lambda & 3 \\ 4 & -4 & 5 - \lambda \end{vmatrix} \\ &= -(5 - \lambda)^2(2 + \lambda) - 3(4)^2 - 4(3)^2 \\ &\quad + 3(4)(2 + \lambda) + 3(4)(5 - \lambda) + 3(4)(5 - \lambda) \\ &= -\lambda^3 + 8\lambda^2 - 17\lambda + 10 \\ &= -(\lambda - 5)(\lambda - 2)(\lambda - 1) = 0, \end{aligned}$$

so the three eigenvalues of A are 1, 2, and 5. To find an eigenvector of A corresponding to the eigenvalue $\lambda = 5$, we must solve the equation $A\mathbf{x} = 5\mathbf{x}$ for \mathbf{x} , which yields the system of equations

$$5x_1 - 3x_2 + 3x_3 = 5x_1,$$

$$4x_1 - 2x_2 + 3x_3 = 5x_2,$$

$$4x_1 - 4x_2 + 5x_3 = 5x_3.$$

The first and third equations imply that $x_2 = x_3$ and $x_1 = x_2$, which, when used in the second equation, yields the identity $x_2 = x_2$. Thus, x_2 is arbitrary, and so any \mathbf{x} having $x_1 = x_2 = x_3$, such as the vector $(1, 1, 1)'$, is an eigenvector of A associated with the eigenvalue 5. In a similar fashion, by solving the equation $A\mathbf{x} = \lambda\mathbf{x}$, when $\lambda = 2$ and $\lambda = 1$, we find that $(1, 1, 0)'$ is an eigenvector corresponding to the eigenvalue 2, and $(0, 1, 1)'$ is an eigenvector corresponding to the eigenvalue 1.

Note that if a nonnull vector \mathbf{x} satisfies (3.1) for a given value of λ , then so will $\alpha\mathbf{x}$ for any nonzero scalar α . Thus, eigenvectors are not uniquely defined unless we impose some scale constraint; for instance, we might only consider eigenvectors, \mathbf{x} , which satisfy $\mathbf{x}'\mathbf{x} = 1$. In this case, for the previous example, we would obtain the three normalized eigenvectors, $(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})'$, $(1/\sqrt{2}, 1/\sqrt{2}, 0)'$, and $(0, 1/\sqrt{2}, 1/\sqrt{2})'$ corresponding to the eigenvalues 5, 2, and 1, respectively. These normalized eigenvectors are unique except for sign, because each of these eigenvectors, when multiplied by -1 , yields another normalized eigenvector.

Example 3.2 illustrates the fact that a real matrix may have complex eigenvalues and eigenvectors.

Example 3.2 The matrix

$$A = \begin{bmatrix} 1 & 1 \\ -2 & -1 \end{bmatrix}$$

has the characteristic equation

$$\begin{aligned} |A - \lambda I_2| &= \begin{vmatrix} 1 - \lambda & 1 \\ -2 & -1 - \lambda \end{vmatrix} \\ &= -(1 - \lambda)(1 + \lambda) + 2 \\ &= \lambda^2 + 1 = 0, \end{aligned}$$

so that the eigenvalues of A are $i = \sqrt{-1}$ and $-i$. To find an eigenvector corresponding to the root i , write $\mathbf{x} = (x_1, x_2)' = (y_1 + iz_1, y_2 + iz_2)'$ and solve for y_1, y_2, z_1, z_2 using the equation $A\mathbf{x} = i\mathbf{x}$. As a result, we find that for any real scalar $\alpha \neq 0$, $\mathbf{x} = (\alpha + i\alpha, -2\alpha)'$ is an eigenvector corresponding to the eigenvalue i . In a similar manner, it can be shown that an eigenvector associated with the eigenvalue $-i$ has the form $\mathbf{x} = (\alpha - i\alpha, -2\alpha)'$.

The m eigenvalues of a matrix A need not all be different because the characteristic equation may have repeated roots. An eigenvalue that occurs as a single solution to the characteristic equation will be called a simple or distinct eigenvalue. Otherwise, an eigenvalue will be called a multiple eigenvalue, and its multiplicity, or more precisely its algebraic multiplicity, will be given by the number of times this solution is repeated.

In some situations, we will find it useful to work with the set of all eigenvectors associated with a specific eigenvalue. This collection, $S_A(\lambda)$, of all eigenvectors corresponding to the particular eigenvalue λ , along with the trivial vector $\mathbf{0}$, is called the eigenspace of A associated with λ ; that is, $S_A(\lambda)$ is given by $S_A(\lambda) = \{\mathbf{x} : \mathbf{x} \in R^m \text{ and } A\mathbf{x} = \lambda\mathbf{x}\}$.

Theorem 3.1 If $S_A(\lambda)$ is the eigenspace of the $m \times m$ matrix A corresponding to the eigenvalue λ , then $S_A(\lambda)$ is a vector subspace of R^m .

Proof. By definition, if $\mathbf{x} \in S_A(\lambda)$, then $A\mathbf{x} = \lambda\mathbf{x}$. Thus, if $\mathbf{x} \in S_A(\lambda)$ and $\mathbf{y} \in S_A(\lambda)$, we have for any scalars α and β

$$A(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha A\mathbf{x} + \beta A\mathbf{y} = \alpha(\lambda\mathbf{x}) + \beta(\lambda\mathbf{y}) = \lambda(\alpha\mathbf{x} + \beta\mathbf{y}).$$

Consequently, $(\alpha\mathbf{x} + \beta\mathbf{y}) \in S_A(\lambda)$, and so $S_A(\lambda)$ is a vector space. \square

Example 3.3 The matrix

$$A = \begin{bmatrix} 2 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

has the characteristic equation

$$\begin{vmatrix} 2 - \lambda & -1 & 0 \\ 0 & 1 - \lambda & 0 \\ 0 & 0 & 1 - \lambda \end{vmatrix} = (1 - \lambda)^2(2 - \lambda) = 0,$$

and so the eigenvalues of A are 1, with multiplicity 2, and 2. To find $S_A(1)$, the eigenspace corresponding to the eigenvalue 1, we solve the equation $A\mathbf{x} = \mathbf{x}$ for \mathbf{x} . We leave it to the reader to verify that two linearly independent solutions result; any solution to $A\mathbf{x} = \mathbf{x}$ will be a linear combination of the two vectors $\mathbf{x}_1 = (0, 0, 1)'$ and $\mathbf{x}_2 = (1, 1, 0)'$. Thus, $S_A(1)$ is the subspace spanned by the basis $\{\mathbf{x}_1, \mathbf{x}_2\}$; that is, $S_A(1)$ is a plane in R^3 . In a similar fashion, we may find the eigenspace, $S_A(2)$. Solving $A\mathbf{x} = 2\mathbf{x}$, we find that \mathbf{x} must be a scalar multiple of $(1, 0, 0)'$. Thus, $S_A(2)$ is the line in R^3 given by $\{(a, 0, 0)' : -\infty < a < \infty\}$.

In Example 3.3, for each of the two values of λ , we have $\dim\{S(\lambda)\}$, which is sometimes referred to as the geometric multiplicity of λ , being equal to the corresponding algebraic multiplicity of λ . This is not always the case. In general, when

we simply use the term *multiplicity*, we will be referring to the algebraic multiplicity. Example 3.4 illustrates that it is possible for $\dim\{S(\lambda)\}$ to be less than the multiplicity of the eigenvalue λ .

Example 3.4 Consider the 3×3 matrix given by

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

Since $|A - \lambda I_3| = (1 - \lambda)^3$, A has the eigenvalue 1 repeated three times. The eigenvalue-eigenvector equation $A\mathbf{x} = \lambda\mathbf{x}$ yields the three scalar equations

$$x_1 + 2x_2 + 3x_3 = x_1,$$

$$x_2 = x_2,$$

$$2x_2 + x_3 = x_3,$$

which have as a solution only vectors of the form $\mathbf{x} = (a, 0, 0)'$. Thus, although the multiplicity of the eigenvalue 1 is 3, the associated eigenspace $S_A(1) = \{(a, 0, 0)' : -\infty < a < \infty\}$ is only one-dimensional.

3.3 SOME BASIC PROPERTIES OF EIGENVALUES AND EIGENVECTORS

In this section, we establish some useful results regarding eigenvalues. The proofs of the results in Theorem 3.2, which are left to the reader as an exercise, are easily obtained by using the characteristic equation or the eigenvalue-eigenvector equation.

Theorem 3.2 Let A be an $m \times m$ matrix. Then the following properties hold:

- (a) The eigenvalues of A' are the same as the eigenvalues of A .
- (b) A is singular if and only if at least one eigenvalue of A is equal to 0.
- (c) The diagonal elements of A are the eigenvalues of A , if A is a triangular matrix.
- (d) The eigenvalues of BAB^{-1} are the same as the eigenvalues of A , if B is a nonsingular $m \times m$ matrix.
- (e) The modulus of each eigenvalue of A is equal to 1 if A is an orthogonal matrix.

We saw in Example 3.4 that it is possible for the dimension of an eigenspace associated with an eigenvalue λ to be less than the multiplicity of λ . Theorem 3.3 shows that if $\dim\{S_A(\lambda)\} \neq r$, where r denotes the multiplicity of λ , then $\dim\{S_A(\lambda)\} < r$.

Theorem 3.3 Suppose that λ is an eigenvalue, with multiplicity $r \geq 1$, of the $m \times m$ matrix A . Then

$$1 \leq \dim\{S_A(\lambda)\} \leq r.$$

Proof. If λ is an eigenvalue of A , then by definition an $x \neq 0$ satisfying the eigenvalue-eigenvector equation $Ax = \lambda x$ exists, and so clearly $\dim\{S_A(\lambda)\} \geq 1$. Now let $k = \dim\{S_A(\lambda)\}$, and let x_1, \dots, x_k be linearly independent eigenvectors corresponding to λ . Form a nonsingular $m \times m$ matrix X that has these k vectors as its first k columns; that is, X has the form $X = [X_1 \ X_2]$, where $X_1 = (x_1, \dots, x_k)$ and X_2 is $m \times (m - k)$. Since each column of X_1 is an eigenvector of A corresponding to the eigenvalue λ , it follows that $AX_1 = \lambda X_1$, and

$$X^{-1}X_1 = \begin{bmatrix} I_k \\ (0) \end{bmatrix}$$

follows from the fact that $X^{-1}X = I_m$. As a result, we find that

$$\begin{aligned} X^{-1}AX &= X^{-1}[AX_1 \ AX_2] \\ &= X^{-1}[\lambda X_1 \ AX_2] \\ &= \begin{bmatrix} \lambda I_k & B_1 \\ (0) & B_2 \end{bmatrix}, \end{aligned}$$

where B_1 and B_2 represent a partitioning of the matrix $X^{-1}AX_2$. If μ is an eigenvalue of $X^{-1}AX$, then

$$\begin{aligned} 0 &= |X^{-1}AX - \mu I_m| = \begin{vmatrix} (\lambda - \mu)I_k & B_1 \\ (0) & B_2 - \mu I_{m-k} \end{vmatrix} \\ &= (\lambda - \mu)^k |B_2 - \mu I_{m-k}|, \end{aligned}$$

where the last equality can be obtained by repeated use of the cofactor expansion formula for a determinant. Thus, λ must be an eigenvalue of $X^{-1}AX$ with multiplicity of at least k . The result now follows because, from Theorem 3.2(d), the eigenvalues of $X^{-1}AX$ are the same as those of A . \square

We now prove Theorem 3.4, which involves both the eigenvalues and the eigenvectors of a matrix.

Theorem 3.4 Let λ be an eigenvalue of the $m \times m$ matrix A , and let x be a corresponding eigenvector. Then,

- (a) if $n \geq 1$ is an integer, λ^n is an eigenvalue of A^n corresponding to the eigenvector x ,

- (b) if A is nonsingular, λ^{-1} is an eigenvalue of A^{-1} corresponding to the eigenvector \mathbf{x} .

Proof. Part (a) can be proven by induction. Clearly (a) holds when $n = 1$ because it follows from the definition of λ and \mathbf{x} . Note that if (a) holds for $n - 1$, that is, $A^{n-1}\mathbf{x} = \lambda^{n-1}\mathbf{x}$, then

$$\begin{aligned} A^n\mathbf{x} &= A(A^{n-1}\mathbf{x}) = A(\lambda^{n-1}\mathbf{x}) \\ &= \lambda^{n-1}(A\mathbf{x}) = \lambda^{n-1}(\lambda\mathbf{x}) = \lambda^n\mathbf{x}, \end{aligned}$$

and so the induction proof of (a) is complete. To prove part (b), premultiply the eigenvalue-eigenvector equation

$$A\mathbf{x} = \lambda\mathbf{x}$$

by A^{-1} , which yields the equation

$$\mathbf{x} = \lambda A^{-1}\mathbf{x}. \quad (3.3)$$

Since A is nonsingular, we know from Theorem 3.2(b) that $\lambda \neq 0$, and so dividing both sides of (3.3) by λ yields

$$A^{-1}\mathbf{x} = \lambda^{-1}\mathbf{x},$$

which is the eigenvalue-eigenvector equation for A^{-1} , with eigenvalue λ^{-1} and eigenvector \mathbf{x} . \square

The determinant and trace of a matrix have simple and useful relationships with the eigenvalues of that matrix. These relationships are established in Theorem 3.5.

Theorem 3.5 Let A be an $m \times m$ matrix with eigenvalues $\lambda_1, \dots, \lambda_m$. Then

- (a) $\text{tr}(A) = \sum_{i=1}^m \lambda_i$,
- (b) $|A| = \prod_{i=1}^m \lambda_i$.

Proof. Recall that the characteristic equation, $|A - \lambda I_m| = 0$, can be expressed in the polynomial form

$$(-\lambda)^m + \alpha_{m-1}(-\lambda)^{m-1} + \dots + \alpha_1(-\lambda) + \alpha_0 = 0. \quad (3.4)$$

We will first identify the coefficients α_0 and α_{m-1} . We can determine α_0 by evaluating the left-hand side of (3.4) at $\lambda = 0$; thus, $\alpha_0 = |A - (0)I_m| = |A|$. To find α_{m-1} , recall that, from its definition, the determinant is expressed as a sum of terms

over all permutations of the integers $(1, 2, \dots, m)$. Since α_{m-1} is the coefficient of $(-\lambda)^{m-1}$, to identify this term we only need to consider the terms in the sum that involve at least $m-1$ of the diagonal elements of $(A - \lambda I_m)$. However, each term in the sum is the product of m elements from the matrix $(A - \lambda I_m)$, multiplied by the appropriate sign, with one element chosen from each row and each column of $(A - \lambda I_m)$. Consequently, the only term in the sum involving at least $m-1$ of the diagonal elements of $(A - \lambda I_m)$ is the term that involves the product of all of the diagonal elements. Since this term involves an even permutation, the sign term will equal $+1$, and so α_{m-1} will be the coefficient of $(-\lambda)^{m-1}$ in

$$(a_{11} - \lambda)(a_{22} - \lambda) \cdots (a_{mm} - \lambda),$$

which clearly is $a_{11} + a_{22} + \cdots + a_{mm}$ or simply $\text{tr}(A)$. Now to relate $\alpha_0 = |A|$ and $\alpha_{m-1} = \text{tr}(A)$ to the eigenvalues of A , note that because $\lambda_1, \dots, \lambda_m$ are the roots to the characteristic equation, which is an m th degree polynomial in λ , it follows that

$$(\lambda_1 - \lambda)(\lambda_2 - \lambda) \cdots (\lambda_m - \lambda) = 0.$$

Multiplying out the left-hand side of this equation and then matching corresponding terms with those in (3.4), we find that

$$|A| = \prod_{i=1}^m \lambda_i, \quad \text{tr}(A) = \sum_{i=1}^m \lambda_i,$$

and this completes the proof. \square

The utility of the formulas for the determinant and trace of a matrix in terms of its eigenvalues is illustrated in the proof of Theorem 3.6.

Theorem 3.6 Let A be an $m \times m$ nonsingular symmetric matrix and \mathbf{c} and \mathbf{d} be $m \times 1$ vectors. Then

$$|A + \mathbf{c}\mathbf{d}'| = |A|(1 + \mathbf{d}'A^{-1}\mathbf{c}).$$

Proof. Since A is nonsingular, we have

$$|A + \mathbf{c}\mathbf{d}'| = |A(I_m + A^{-1}\mathbf{c}\mathbf{d}')| = |A||I_m + \mathbf{b}\mathbf{d}'|,$$

where $\mathbf{b} = A^{-1}\mathbf{c}$. For any \mathbf{x} orthogonal to \mathbf{d} , we have

$$(I_m + \mathbf{b}\mathbf{d}')\mathbf{x} = \mathbf{x},$$

so 1 is an eigenvalue of $I_m + \mathbf{b}\mathbf{d}'$ with multiplicity at least $m-1$ because there are $m-1$ linearly independent vectors orthogonal to \mathbf{d} . However, $\text{tr}(I_m + \mathbf{b}\mathbf{d}') = m + \mathbf{d}'\mathbf{b}$, which implies that the final eigenvalue of $I_m + \mathbf{b}\mathbf{d}'$ is given by $1 + \mathbf{d}'\mathbf{b}$. Taking

the product of these eigenvalues, we find that $|I_m + \mathbf{b}\mathbf{d}'| = (1 + \mathbf{d}'\mathbf{b})$. The result now follows because $\mathbf{d}'\mathbf{b} = \mathbf{d}'A^{-1}\mathbf{c}$. \square

Theorem 3.7 gives a sufficient condition for a set of eigenvectors to be linearly independent.

Theorem 3.7 Suppose $\mathbf{x}_1, \dots, \mathbf{x}_r$ are eigenvectors of the $m \times m$ matrix A , where $r \leq m$. If the corresponding eigenvalues $\lambda_1, \dots, \lambda_r$ are such that $\lambda_i \neq \lambda_j$ for all $i \neq j$, then the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ are linearly independent.

Proof. Our proof is by contradiction; that is, we begin by assuming that the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ are linearly dependent. Let h be the largest integer for which $\mathbf{x}_1, \dots, \mathbf{x}_h$ are linearly independent. Such a set can be found because \mathbf{x}_1 , being an eigenvector, cannot equal $\mathbf{0}$, and so it is linearly independent. The vectors $\mathbf{x}_1, \dots, \mathbf{x}_{h+1}$ must be linearly dependent, so there exist scalars $\alpha_1, \dots, \alpha_{h+1}$, with at least two not equal to zero because no eigenvector can be the null vector, such that

$$\alpha_1 \mathbf{x}_1 + \dots + \alpha_{h+1} \mathbf{x}_{h+1} = \mathbf{0}.$$

Premultiplying the left-hand side of this equation by $(A - \lambda_{h+1}I_m)$, we find that

$$\begin{aligned} & \alpha_1(A - \lambda_{h+1}I_m)\mathbf{x}_1 + \dots + \alpha_{h+1}(A - \lambda_{h+1}I_m)\mathbf{x}_{h+1} \\ &= \alpha_1(A\mathbf{x}_1 - \lambda_{h+1}\mathbf{x}_1) + \dots + \alpha_{h+1}(A\mathbf{x}_{h+1} - \lambda_{h+1}\mathbf{x}_{h+1}) \\ &= \alpha_1(\lambda_1 - \lambda_{h+1})\mathbf{x}_1 + \dots + \alpha_h(\lambda_h - \lambda_{h+1})\mathbf{x}_h \end{aligned}$$

also must be equal to $\mathbf{0}$. But $\mathbf{x}_1, \dots, \mathbf{x}_h$ are linearly independent, so it follows that

$$\alpha_1(\lambda_1 - \lambda_{h+1}) = \dots = \alpha_h(\lambda_h - \lambda_{h+1}) = 0.$$

We know that at least one of the scalars $\alpha_1, \dots, \alpha_h$ is not equal to zero, and if, for instance, α_i is one of these nonzero scalars, we then must have $\lambda_i = \lambda_{h+1}$. This result contradicts the conditions of the theorem, so the vectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ must be linearly independent. \square

If the eigenvalues $\lambda_1, \dots, \lambda_m$ of an $m \times m$ matrix A are all distinct, then it follows from Theorem 3.7 that the matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, where \mathbf{x}_i is an eigenvector corresponding to λ_i , is nonsingular. It also follows from the eigenvalue-eigenvector equation $A\mathbf{x}_i = \lambda_i\mathbf{x}_i$ that if we define the diagonal matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$, then $AX = X\Lambda$. Premultiplying this equation by X^{-1} yields the identity $X^{-1}AX = \Lambda$. Any square matrix that can be transformed into a diagonal matrix through the postmultiplication by a nonsingular matrix and premultiplication by its inverse is said to be diagonalizable. Thus, a square matrix with distinct eigenvalues is diagonalizable.

When X is nonsingular, the equation $AX = X\Lambda$ can also be rearranged as $A = X\Lambda X^{-1}$. That is, in this case, A can be determined from its eigenvalues and any set of corresponding linearly independent eigenvectors.

Example 3.5 Suppose that A is a 2×2 matrix with eigenvalues 1 and 2, and corresponding eigenvectors given as the columns of the matrix

$$X = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}.$$

Since $|X| = 1$, X is nonsingular and

$$X^{-1} = \begin{bmatrix} 2 & -3 \\ -3 & 5 \end{bmatrix}.$$

Thus, we have enough information to compute A ; that is,

$$\begin{aligned} A &= X\Lambda X^{-1} = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix} \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix} \begin{bmatrix} 2 & -3 \\ -3 & 5 \end{bmatrix} \\ &= \begin{bmatrix} -8 & 15 \\ -6 & 11 \end{bmatrix}. \end{aligned}$$

Consider a second 2×2 matrix B that has both of its eigenvalues equal to 0, but only a single linearly independent eigenvector being any nonzero scalar multiple of $e_1 = (1, 0)'$. The eigenvalue-eigenvector equation $Be_1 = \mathbf{0}$ implies that $b_{11} = b_{21} = 0$, and so the characteristic equation for B is

$$\begin{vmatrix} -\lambda & b_{12} \\ 0 & b_{22} - \lambda \end{vmatrix} = -\lambda(b_{22} - \lambda) = 0.$$

Since both of the eigenvalues of B are 0, we must have $b_{22} = 0$, and so B is of the form

$$B = \begin{bmatrix} 0 & b_{12} \\ 0 & 0 \end{bmatrix},$$

where $b_{12} \neq 0$, because otherwise B would have two linearly independent eigenvectors. Note, however, that the value of b_{12} cannot be determined.

Clearly, when a matrix is diagonalizable, then its rank equals the number of its nonzero eigenvalues, because

$$\text{rank}(A) = \text{rank}(X^{-1}AX) = \text{rank}(\Lambda)$$

follows from Theorem 1.10. This relationship between the number of nonzero eigenvalues and the rank of a square matrix does not necessarily hold if the matrix is not diagonalizable.

Example 3.6 Consider the 2×2 matrices

$$A = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix} \quad \text{and} \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

Clearly, both A and B have rank of 1. Now the characteristic equation of A simplifies to $\lambda(1 - \lambda) = 0$ so that the eigenvalues of A are 0 and 1, and thus, in this case, $\text{rank}(A)$ equals the number of nonzero eigenvalues. The characteristic equation for B simplifies to $\lambda^2 = 0$, so B has the eigenvalue 0 repeated twice. Hence, the rank of B exceeds its number of nonzero eigenvalues.

Theorem 3.8, known as the Cayley–Hamilton theorem, states that a matrix satisfies its own characteristic equation.

Theorem 3.8 Let A be an $m \times m$ matrix with eigenvalues $\lambda_1, \dots, \lambda_m$. Then

$$\prod_{i=1}^m (A - \lambda_i I_m) = (0);$$

that is, if $(-\lambda)^m + \alpha_{m-1}(-\lambda)^{m-1} + \dots + \alpha_1(-\lambda) + \alpha_0 = 0$ is the characteristic equation of A , then

$$(-A)^m + \alpha_{m-1}(-A)^{m-1} + \dots + \alpha_1(-A) + \alpha_0 I_m = (0).$$

Proof. If $(A - \lambda I_m)_\#$ is the adjoint of $A - \lambda I_m$, then

$$(A - \lambda I_m)(A - \lambda I_m)_\# = |A - \lambda I_m| I_m.$$

Since $|A - \lambda I_m|$ is an m th degree polynomial in λ , it follows that $(A - \lambda I_m)_\#$ is a polynomial in λ of degree at most $m - 1$. That is, we can write

$$(A - \lambda I_m)_\# = \sum_{i=0}^{m-1} (-\lambda)^i B_i,$$

where B_i is an $m \times m$ matrix for $i = 0, \dots, m - 1$, and so

$$\begin{aligned} (A - \lambda I_m)(A - \lambda I_m)_\# &= (A - \lambda I_m) \sum_{i=0}^{m-1} (-\lambda)^i B_i \\ &= AB_0 + \sum_{i=1}^{m-1} (-\lambda)^i (AB_i + B_{i-1}) + (-\lambda)^m B_{m-1}. \end{aligned}$$

Equating this to $|A - \lambda I_m|I_m = \sum_{i=0}^{m-1} (-\lambda)^i \alpha_i I_m + (-\lambda)^m I_m$, we find that

$$\begin{aligned} AB_0 &= \alpha_0 I_m, \\ AB_1 + B_0 &= \alpha_1 I_m, \\ &\vdots \\ AB_{m-1} + B_{m-2} &= \alpha_{m-1} I_m, \\ B_{m-1} &= I_m. \end{aligned}$$

Adding these equations after multiplying the i th equation by $(-A)^{i-1}$, where $(-A)^0 = I_m$, yields the desired result. \square

3.4 SYMMETRIC MATRICES

Many of the applications involving eigenvalues and eigenvectors in statistics are ones that deal with a symmetric matrix, and symmetric matrices have some especially nice properties regarding eigenvalues and eigenvectors. In this section, we will develop some of these properties.

We have seen that a matrix may have complex eigenvalues even when the matrix itself is real. This is not the case for symmetric matrices.

Theorem 3.9 Let A be an $m \times m$ real symmetric matrix. Then the eigenvalues of A are real, and corresponding to any eigenvalue, eigenvectors that are real exist.

Proof. Let $\lambda = \alpha + i\beta$ be an eigenvalue of A and $\mathbf{x} = \mathbf{y} + i\mathbf{z}$ a corresponding eigenvector, where $i = \sqrt{-1}$. We will first show that $\beta = 0$. Substitution of these expressions for λ and \mathbf{x} in the eigenvalue-eigenvector equation $A\mathbf{x} = \lambda\mathbf{x}$ yields

$$A(\mathbf{y} + i\mathbf{z}) = (\alpha + i\beta)(\mathbf{y} + i\mathbf{z}). \quad (3.5)$$

Premultiplying (3.5) by $(\mathbf{y} - i\mathbf{z})'$, we get

$$(\mathbf{y} - i\mathbf{z})' A(\mathbf{y} + i\mathbf{z}) = (\alpha + i\beta)(\mathbf{y} - i\mathbf{z})'(\mathbf{y} + i\mathbf{z}),$$

which simplifies to

$$\mathbf{y}' A \mathbf{y} + \mathbf{z}' A \mathbf{z} = (\alpha + i\beta)(\mathbf{y}' \mathbf{y} + \mathbf{z}' \mathbf{z}),$$

because $\mathbf{y}' A \mathbf{z} = \mathbf{z}' A \mathbf{y}$ follows from the symmetry of A . Now $\mathbf{x} \neq \mathbf{0}$ implies that $(\mathbf{y}' \mathbf{y} + \mathbf{z}' \mathbf{z}) > 0$ and, consequently, we must have $\beta = 0$ because the left-hand side of the equation above is real. Substituting $\beta = 0$ in (3.5), we find that

$$A\mathbf{y} + iA\mathbf{z} = \alpha\mathbf{y} + i\alpha\mathbf{z}.$$

Thus, $\mathbf{x} = \mathbf{y} + iz$ will be an eigenvector of A corresponding to $\lambda = \alpha$ as long as \mathbf{y} and \mathbf{z} satisfy $A\mathbf{y} = \alpha\mathbf{y}$, $A\mathbf{z} = \alpha\mathbf{z}$, and at least one is not $\mathbf{0}$, so that $\mathbf{x} \neq \mathbf{0}$. A real eigenvector is then constructed by selecting $\mathbf{y} \neq \mathbf{0}$, such that $A\mathbf{y} = \alpha\mathbf{y}$, and $\mathbf{z} = \mathbf{0}$. \square

We have seen that a set of eigenvectors of an $m \times m$ matrix A is linearly independent if the associated eigenvalues are all different from one other. We will now show that, if A is symmetric, we can say a bit more. Suppose that \mathbf{x} and \mathbf{y} are eigenvectors of A corresponding to the eigenvalues λ and γ , where $\lambda \neq \gamma$. Then, because A is symmetric, it follows that

$$\begin{aligned}\lambda \mathbf{x}'\mathbf{y} &= (\lambda \mathbf{x})'\mathbf{y} = (A\mathbf{x})'\mathbf{y} = \mathbf{x}'A'\mathbf{y} \\ &= \mathbf{x}'(A\mathbf{y}) = \mathbf{x}'(\gamma\mathbf{y}) = \gamma \mathbf{x}'\mathbf{y}.\end{aligned}$$

Since $\lambda \neq \gamma$, we must have $\mathbf{x}'\mathbf{y} = 0$; that is, eigenvectors corresponding to different eigenvalues must be orthogonal. Thus, if the m eigenvalues of A are distinct, then the set of corresponding eigenvectors will form a group of mutually orthogonal vectors. We will show that this is still possible when A has multiple eigenvalues. Before we prove this, we will need Theorem 3.10.

Theorem 3.10 Let A be an $m \times m$ symmetric matrix, and let \mathbf{x} be any nonzero $m \times 1$ vector. Then for some $r \geq 1$, the vector space spanned by the vectors $\mathbf{x}, A\mathbf{x}, \dots, A^{r-1}\mathbf{x}$ contains an eigenvector of A .

Proof. Let r be the smallest integer for which $\mathbf{x}, A\mathbf{x}, \dots, A^r\mathbf{x}$ form a linearly dependent set. Then scalars $\alpha_0, \dots, \alpha_r$ exist, not all of which are zero, such that

$$\alpha_0\mathbf{x} + \alpha_1A\mathbf{x} + \dots + \alpha_rA^r\mathbf{x} = (\alpha_0I_m + \alpha_1A + \dots + A^r)\mathbf{x} = \mathbf{0},$$

where, without loss of generality, we have taken $\alpha_r = 1$, because the way r was chosen guarantees that α_r is not zero. The expression in the parentheses is an r th-degree matrix polynomial in A , which can be factored in a fashion similar to the way scalar polynomials are factored; that is, it can be written as

$$(A - \gamma_1I_m)(A - \gamma_2I_m) \cdots (A - \gamma_rI_m),$$

where $\gamma_1, \dots, \gamma_r$ are the roots of the polynomial which satisfy $\alpha_0 = (-1)^r\gamma_1\gamma_2 \cdots \gamma_r$, \dots , $\alpha_{r-1} = -(\gamma_1 + \gamma_2 + \dots + \gamma_r)$. Let

$$\begin{aligned}\mathbf{y} &= (A - \gamma_2I_m) \cdots (A - \gamma_rI_m)\mathbf{x} \\ &= (-1)^{r-1}\gamma_2 \cdots \gamma_r\mathbf{x} + \dots + A^{r-1}\mathbf{x},\end{aligned}$$

and note that $\mathbf{y} \neq \mathbf{0}$ because, otherwise, $\mathbf{x}, A\mathbf{x}, \dots, A^{r-1}\mathbf{x}$ would be a linearly dependent set, contradicting the definition of r . Thus, \mathbf{y} is in the space spanned by $\mathbf{x}, A\mathbf{x}, \dots, A^{r-1}\mathbf{x}$ and

$$(A - \gamma_1I_m)\mathbf{y} = (A - \gamma_1I_m)(A - \gamma_2I_m) \cdots (A - \gamma_rI_m)\mathbf{x} = \mathbf{0}.$$

Consequently, \mathbf{y} is an eigenvector of A corresponding to the eigenvalue γ_1 , and so the proof is complete. \square

Theorem 3.11 If the $m \times m$ matrix A is symmetric, then it is possible to construct a set of m eigenvectors of A such that the set is orthonormal.

Proof. We first show that if we have an orthonormal set of eigenvectors, say $\mathbf{x}_1, \dots, \mathbf{x}_h$, where $1 \leq h < m$, then we can find another normalized eigenvector \mathbf{x}_{h+1} orthogonal to each of these vectors. Select any vector \mathbf{x} that is orthogonal to each of the vectors $\mathbf{x}_1, \dots, \mathbf{x}_h$. Note that for any positive integer k , $A^k \mathbf{x}$ is also orthogonal to $\mathbf{x}_1, \dots, \mathbf{x}_h$ because, if λ_i is the eigenvalue corresponding to \mathbf{x}_i , it follows from the symmetry of A and Theorem 3.4(a) that

$$\mathbf{x}'_i A^k \mathbf{x} = \{(A^k)' \mathbf{x}_i\}' \mathbf{x} = (A^k \mathbf{x}_i)' \mathbf{x} = \lambda_i^k \mathbf{x}'_i \mathbf{x} = 0.$$

From the previous theorem, we know that, for some r , the space spanned by the vectors $\mathbf{x}, A\mathbf{x}, \dots, A^{r-1}\mathbf{x}$ contains an eigenvector, say \mathbf{y} , of A . This vector \mathbf{y} also must be orthogonal to $\mathbf{x}_1, \dots, \mathbf{x}_h$ because it is from a vector space spanned by a set of vectors orthogonal to $\mathbf{x}_1, \dots, \mathbf{x}_h$. Thus, we can take $\mathbf{x}_{h+1} = (\mathbf{y}'\mathbf{y})^{-1/2}\mathbf{y}$. The theorem now follows by starting with any eigenvector of A , and then using the previous argument $m - 1$ times. \square

If we let the $m \times m$ matrix $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, where $\mathbf{x}_1, \dots, \mathbf{x}_m$ are the orthonormal vectors described in the proof, and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$, then the eigenvalue-eigenvector equations $A\mathbf{x}_i = \lambda_i \mathbf{x}_i$ for $i = 1, \dots, m$ can be expressed collectively as the matrix equation $AX = X\Lambda$. Since the columns of X are orthonormal vectors, X is an orthogonal matrix. Premultiplication of our matrix equation by X' yields the relationship $X'AX = \Lambda$, or equivalently

$$A = X\Lambda X',$$

which is known as the spectral decomposition of A . We will see in Section 4.2 that there is a very useful generalization of this decomposition, known as the singular value decomposition, which holds for any $m \times n$ matrix A ; in particular, there exist $m \times m$ and $n \times n$ orthogonal matrices P and Q and an $m \times n$ matrix D with $d_{ij} = 0$ if $i \neq j$, such that $A = PDQ'$.

Note that it follows from Theorem 3.2(d) that the eigenvalues of A are the same as the eigenvalues of Λ , which are the diagonal elements of Λ . Thus, if λ is a multiple eigenvalue of A with multiplicity $r > 1$, then r of the diagonal elements of Λ are equal to λ and r of the eigenvectors, say $\mathbf{x}_1, \dots, \mathbf{x}_r$, correspond to this eigenvalue λ . Consequently, the dimension of the eigenspace of A , $S_A(\lambda)$, corresponding to λ , is equal to the multiplicity r . The set of orthonormal eigenvectors corresponding to this eigenvalue is not unique. Any orthonormal basis for $S_A(\lambda)$ will be a set of r orthonormal vectors associated with the eigenvalue λ . For example, if we let $X_1 =$

$(\mathbf{x}_1, \dots, \mathbf{x}_r)$ and let Q be any $r \times r$ orthogonal matrix, the columns of $Y_1 = X_1 Q$ also form a set of orthonormal eigenvectors corresponding to λ .

Example 3.7 One application of an eigenanalysis in statistics involves overcoming difficulties associated with a regression analysis in which the explanatory variables are nearly linearly dependent. This situation is often referred to as multicollinearity. In this case, some of the explanatory variables are providing redundant information about the response variable. As a result, the least squares estimator of β in the model $\mathbf{y} = X\beta + \epsilon$

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

will be imprecise because its covariance matrix

$$\begin{aligned}\text{var}(\hat{\beta}) &= (X'X)^{-1}X'\{\text{var}(\mathbf{y})\}X(X'X)^{-1} \\ &= (X'X)^{-1}X'\{\sigma^2 I\}X(X'X)^{-1} \\ &= \sigma^2(X'X)^{-1}\end{aligned}$$

will tend to have some large elements because of the near singularity of $X'X$. If the near linear dependence is simply because one of the explanatory variables, say x_j , is nearly a scalar multiple of another, say x_l , then one could simply eliminate one of these variables from the model. However, in most cases, the near linear dependence is not this straightforward. We will see that an eigenanalysis will help reveal any of these dependencies. Suppose that we standardize the explanatory variables so that we have the model

$$\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \delta_1 + \epsilon$$

discussed in Example 2.16. Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_k)$ contain the eigenvalues of $Z_1'Z_1$ in descending order of magnitude, and let U be an orthogonal matrix that has corresponding normalized eigenvectors of $Z_1'Z_1$ as its columns, so that $Z_1'Z_1 = U\Lambda U'$. It was shown in Example 2.16 that the estimation of \mathbf{y} is unaffected by a non-singular transformation of the explanatory variables; that is, we could just as well work with the model

$$\mathbf{y} = \alpha_0 \mathbf{1}_N + W_1 \alpha_1 + \epsilon,$$

where $\alpha_0 = \delta_0$, $\alpha_1 = T^{-1}\delta_1$, $W_1 = Z_1 T$, and T is a nonsingular matrix. A method, referred to as principal components regression, deals with the problems associated with multicollinearity by using the orthogonal transformations $W_1 = Z_1 U$ and $\alpha_1 = U'\delta_1$ of the standardized explanatory variables and parameter vector. The k new explanatory variables are called the principal components; the variable corresponding to the i th column of W_1 is called the i th principal component. Since $W_1'W_1 = U'Z_1'Z_1U = \Lambda$ and $\mathbf{1}_N'W_1 = \mathbf{1}_N'Z_1U = \mathbf{0}'U = \mathbf{0}'$, the least squares estimate of α_1 is

$$\hat{\alpha}_1 = (W_1'W_1)^{-1}W_1'\mathbf{y} = \Lambda^{-1}W_1'\mathbf{y},$$

whereas its covariance matrix simplifies to

$$\text{var}(\hat{\alpha}_1) = \sigma^2(W_1'W_1)^{-1} = \sigma^2\Lambda^{-1}.$$

If $Z_1'Z_1$ and, hence, also $W_1'W_1$ is nearly singular, then at least one of the λ_i 's will be very small, whereas the variances of the corresponding $\hat{\alpha}_i$'s will be very large. Since the explanatory variables have been standardized, $W_1'W_1$ is $N - 1$ times the sample correlation matrix of the principal components computed from the N observations. Thus, if $\lambda_i \approx 0$, then the i th principal component is nearly constant from observation to observation, and so it contributes little to the estimation of \mathbf{y} . If $\lambda_i \approx 0$ for $i = k - r + 1, \dots, k$, then the problems associated with multicollinearity can be avoided by eliminating the last r principal components from the model; in other words, the principal components regression model is

$$\mathbf{y} = \alpha_0 \mathbf{1}_N + W_{11} \alpha_{11} + \epsilon,$$

where W_{11} and α_{11}' are obtained from W_1 and α_1' by deleting their last r columns. If we let $\Lambda_1 = \text{diag}(\lambda_1, \dots, \lambda_{k-r})$, then the least squares estimate of α_{11} can be written as

$$\hat{\alpha}_{11} = (W_{11}'W_{11})^{-1}W_{11}'\mathbf{y} = \Lambda_1^{-1}W_{11}'\mathbf{y}.$$

Note that because of the orthogonality of the principal components, $\hat{\alpha}_{11}$ is identical to the first $k - r$ components of $\hat{\alpha}_1$. The estimate $\hat{\alpha}_{11}$ can be used to find the principal components estimate of δ_1 in the original standardized model. Recall that δ_1 and α_1 are related through the identity $\delta_1 = U\alpha_1$. By eliminating the last r principal components, we are replacing this identity with the identity $\delta_1 = U_1\alpha_{11}$, where $U = [U_1 \ U_2]$ and U_1 is $k \times (k - r)$. Thus, the principal components regression estimate of δ_1 is given by

$$\hat{\delta}_{1*} = U_1 \hat{\alpha}_{11} = U_1 \Lambda_1^{-1} W_{11}' \mathbf{y}.$$

A set of orthonormal eigenvectors of a matrix A can be used to find what are known as the eigenprojections of A .

Definition 3.1 Let λ be an eigenvalue of the $m \times m$ symmetric matrix A with multiplicity $r \geq 1$. If $\mathbf{x}_1, \dots, \mathbf{x}_r$ is a set of orthonormal eigenvectors corresponding to λ , then the eigenprojection of A associated with the eigenvalue λ is given by

$$P_A(\lambda) = \sum_{i=1}^r \mathbf{x}_i \mathbf{x}_i'.$$

The eigenprojection $P_A(\lambda)$ is simply the projection matrix for the vector space $S_A(\lambda)$. Thus, for any $\mathbf{x} \in R^m$, $\mathbf{y} = P_A(\lambda)\mathbf{x}$ gives the orthogonal projection of \mathbf{x} onto the eigenspace $S_A(\lambda)$. If we define X_1 as before, that is $X_1 = (\mathbf{x}_1, \dots, \mathbf{x}_r)$,

then $P_A(\lambda) = X_1 X_1'$. Note that $P_A(\lambda)$ is unique even though the set of eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_r$ is not unique; for instance, if $Y_1 = X_1 Q$, where Q is an arbitrary $r \times r$ orthogonal matrix, then the columns of Y_1 form another set of orthonormal eigenvectors corresponding to λ , but

$$\begin{aligned} Y_1 Y_1' &= (X_1 Q)(X_1 Q)' = X_1 Q Q' X_1' \\ &= X_1 I_r X_1' = X_1 X_1' = P_A(\lambda). \end{aligned}$$

The term *spectral decomposition* comes from the term *spectral set of A* for the set of all eigenvalues of A excluding repetitions of the same value. Suppose the $m \times m$ matrix A has the spectral set $\{\mu_1, \dots, \mu_k\}$, where $k \leq m$, because some of the μ_i 's may correspond to multiple eigenvalues. The set of μ_i 's may be different from our set of λ_i 's in that we do not repeat values for the μ_i 's. Thus, if A is 4×4 with eigenvalues $\lambda_1 = 3, \lambda_2 = 2, \lambda_3 = 2$, and $\lambda_4 = 1$, then the spectral set of A is $\{3, 2, 1\}$. Using X and Λ as previously defined, the spectral decomposition states that

$$A = X \Lambda X' = \sum_{i=1}^m \lambda_i \mathbf{x}_i \mathbf{x}_i' = \sum_{i=1}^k \mu_i P_A(\mu_i),$$

so that A has been decomposed into a sum of terms, one corresponding to each value in the spectral set. If m_i is the multiplicity of μ_i and $\lambda_{M_i+1} = \dots = \lambda_{M_i+m_i} = \mu_i$, where $M_1 = 0$ and $M_i = \sum_{j=1}^{i-1} m_j$ for $i = 2, \dots, k$, then we have $P_A(\mu_i) = \sum_{j=1}^{m_i} \mathbf{x}_{M_i+j} \mathbf{x}_{M_i+j}'$. Note that when we write the decomposition as $A = \sum_{i=1}^k \mu_i P_A(\mu_i)$, the terms in the sum are uniquely defined because of the uniqueness of the projection matrix of a vector space. On the other hand, the decomposition $A = \sum_{i=1}^m \lambda_i \mathbf{x}_i \mathbf{x}_i'$ does not have uniquely defined terms unless the λ_i 's are distinct.

Example 3.8 It can be easily verified by solving the characteristic equation for the 3×3 symmetric matrix

$$A = \begin{bmatrix} 5 & -1 & -1 \\ -1 & 5 & -1 \\ -1 & -1 & 5 \end{bmatrix}$$

that A has the simple eigenvalue 3 and the multiple eigenvalue 6, with multiplicity 2. The unique (except for sign) unit eigenvector associated with the eigenvalue 3 can be shown to equal $(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})'$, whereas a set of orthonormal eigenvectors associated with 6 is given by $(-2/\sqrt{6}, 1/\sqrt{6}, 1/\sqrt{6})'$ and $(0, 1/\sqrt{2}, -1/\sqrt{2})'$. Thus,

the spectral decomposition of A is given by

$$\begin{bmatrix} 5 & -1 & -1 \\ -1 & 5 & -1 \\ -1 & -1 & 5 \end{bmatrix} = \begin{bmatrix} 1/\sqrt{3} & -2/\sqrt{6} & 0 \\ 1/\sqrt{3} & 1/\sqrt{6} & 1/\sqrt{2} \\ 1/\sqrt{3} & 1/\sqrt{6} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \end{bmatrix} \\ \times \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ -2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix},$$

and the two eigenprojections of A are

$$P_A(3) = \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} [1/\sqrt{3} \quad 1/\sqrt{3} \quad 1/\sqrt{3}] = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \\ P_A(6) = \begin{bmatrix} -2/\sqrt{6} & 0 \\ 1/\sqrt{6} & 1/\sqrt{2} \\ 1/\sqrt{6} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} -2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \\ = \frac{1}{3} \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix}.$$

The relationship between the rank of a matrix and the number of its nonzero eigenvalues becomes an exact one for symmetric matrices.

Theorem 3.12 Suppose that the $m \times m$ matrix A has r nonzero eigenvalues. Then, if A is symmetric, $\text{rank}(A) = r$.

Proof. If $A = X\Lambda X'$ is the spectral decomposition of A , then the diagonal matrix Λ has r nonzero diagonal elements and

$$\text{rank}(A) = \text{rank}(X\Lambda X') = \text{rank}(\Lambda),$$

because the multiplication of a matrix by nonsingular matrices does not affect the rank. Clearly, the rank of a diagonal matrix equals the number of its nonzero diagonal elements, so the result follows. \square

Some of the most important applications of eigenvalues and eigenvectors in statistics involve the analysis of covariance and correlation matrices.

Example 3.9 In some situations, a matrix has some special structure that when recognized, can be used to expedite the calculation of eigenvalues and eigenvectors. In this example, we consider a structure sometimes possessed by an $m \times m$ covariance

matrix. This structure is one that has equal variances and equal correlations; that is, the covariance matrix has the form

$$\Omega = \sigma^2 \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix}.$$

Alternatively, Ω can be expressed as $\Omega = \sigma^2\{(1 - \rho)I_m + \rho\mathbf{1}_m\mathbf{1}_m'\}$, so that it is a function of the vector $\mathbf{1}_m$. This vector also plays a crucial role in the eigenanalysis of Ω because

$$\Omega\mathbf{1}_m = \sigma^2\{(1 - \rho)\mathbf{1}_m + \rho\mathbf{1}_m\mathbf{1}_m'\mathbf{1}_m\} = \sigma^2\{(1 - \rho) + m\rho\}\mathbf{1}_m.$$

Thus, $\mathbf{1}_m$ is an eigenvector of Ω corresponding to the eigenvalue $\sigma^2\{(1 - \rho) + m\rho\}$. The remaining eigenvalues of Ω can be identified by noting that if \mathbf{x} is any $m \times 1$ vector orthogonal to $\mathbf{1}_m$, then

$$\Omega\mathbf{x} = \sigma^2\{(1 - \rho)\mathbf{x} + \rho\mathbf{1}_m\mathbf{1}_m'\mathbf{x}\} = \sigma^2(1 - \rho)\mathbf{x},$$

and so \mathbf{x} is an eigenvector of Ω corresponding to the eigenvalue $\sigma^2(1 - \rho)$. Since there are $m - 1$ linearly independent vectors orthogonal to $\mathbf{1}_m$, the eigenvalue $\sigma^2(1 - \rho)$ is repeated $m - 1$ times. The order of these two distinct eigenvalues depends on the value of ρ ; $\sigma^2\{(1 - \rho) + m\rho\}$ will be larger than $\sigma^2(1 - \rho)$ only if ρ is positive.

Example 3.10 A covariance matrix can be any symmetric nonnegative definite matrix. Consequently, for a given set of m nonnegative numbers and a given set of m orthonormal $m \times 1$ vectors, it is possible to construct an $m \times m$ covariance matrix with these numbers and vectors as its eigenvalues and eigenvectors. On the other hand, a correlation matrix has the additional constraint that its diagonal elements must each equal 1, and this extra restriction has an impact on the eigenanalysis of correlation matrices; that is, a much more limited set of possible eigenvalues and eigenvectors exists for correlation matrices. For the most extreme case, consider a 2×2 correlation matrix that must have the form,

$$P = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

with $-1 \leq \rho \leq 1$, because P must be nonnegative definite. The characteristic equation $|P - \lambda I_2| = 0$ readily admits the two eigenvalues $1 + \rho$ and $1 - \rho$. Using these in the eigenvalue-eigenvector equation $P\mathbf{x} = \lambda\mathbf{x}$, we find that regardless of the value of ρ , $(1/\sqrt{2}, 1/\sqrt{2})'$ must be an eigenvector corresponding to $1 + \rho$, whereas $(1/\sqrt{2}, -1/\sqrt{2})'$ must be an eigenvector corresponding to $1 - \rho$. Thus, ignoring sign changes, only one set of orthonormal eigenvectors is possible for a 2×2 correlation matrix if $\rho \neq 0$. This number of possible sets of orthonormal eigenvectors

increases as the order m increases. In some situations, such as simulation studies of analyses of correlation matrices, one may wish to construct a correlation matrix with some particular structure with regard to its eigenvalues or eigenvectors. For example, suppose that we want to construct an $m \times m$ correlation matrix that has three distinct eigenvalues with one of them being repeated $m - 2$ times. Thus, this correlation matrix has the form

$$P = \lambda_1 \mathbf{x}_1 \mathbf{x}_1' + \lambda_2 \mathbf{x}_2 \mathbf{x}_2' + \sum_{i=3}^m \lambda \mathbf{x}_i \mathbf{x}_i',$$

where λ_1 , λ_2 , and λ are the distinct eigenvalues of P , and $\mathbf{x}_1, \dots, \mathbf{x}_m$ are corresponding normalized eigenvectors. Since P is nonnegative definite, we must have $\lambda_1 \geq 0$, $\lambda_2 \geq 0$, and $\lambda \geq 0$, whereas $\text{tr}(P) = m$ implies that $\lambda = (m - \lambda_1 - \lambda_2)/(m - 2)$. Note that P can be written as

$$P = (\lambda_1 - \lambda) \mathbf{x}_1 \mathbf{x}_1' + (\lambda_2 - \lambda) \mathbf{x}_2 \mathbf{x}_2' + \lambda I_m,$$

so that the constraint $(P)_{ii} = 1$ implies that

$$(\lambda_1 - \lambda) x_{i1}^2 + (\lambda_2 - \lambda) x_{i2}^2 + \lambda = 1$$

or, equivalently,

$$x_{i2}^2 = \frac{1 - \lambda - (\lambda_1 - \lambda) x_{i1}^2}{(\lambda_2 - \lambda)}.$$

These constraints can then be used to construct a particular matrix. For instance, suppose that we want to construct a 4×4 correlation matrix with eigenvalues $\lambda_1 = 2$, $\lambda_2 = 1$, and $\lambda = 0.5$ repeated twice. If we choose $\mathbf{x}_1 = (0.5, 0.5, 0.5, 0.5)'$, then we must have $x_{i2}^2 = 0.25$, and so because of the orthogonality of \mathbf{x}_1 and \mathbf{x}_2 , \mathbf{x}_2 can be any vector obtained from \mathbf{x}_1 by negating two of its components. For example, if we take $\mathbf{x}_2 = (0.5, -0.5, 0.5, -0.5)'$, then

$$P = \begin{bmatrix} 1 & 0.25 & 0.50 & 0.25 \\ 0.25 & 1 & 0.25 & 0.50 \\ 0.50 & 0.25 & 1 & 0.25 \\ 0.25 & 0.50 & 0.25 & 1 \end{bmatrix}.$$

3.5 CONTINUITY OF EIGENVALUES AND EIGENPROJECTIONS

Our first result of this section is one that bounds the absolute difference between eigenvalues of two matrices by a function of the absolute differences of the elements of the two matrices. A proof of Theorem 3.13 can be found in Ostrowski (1973). For some other similar bounds, see Elsner (1982).

Theorem 3.13 Let A and B be $m \times m$ matrices with eigenvalues $\lambda_1, \dots, \lambda_m$ and $\gamma_1, \dots, \gamma_m$, respectively. Define

$$M = \max_{1 \leq i \leq m, 1 \leq j \leq m} (|a_{ij}|, |b_{ij}|)$$

and

$$\delta(A, B) = \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^m |a_{ij} - b_{ij}|.$$

Then

$$\max_{1 \leq i \leq m} \min_{1 \leq j \leq m} |\lambda_i - \gamma_j| \leq (m+2)M^{1-1/m} \delta(A, B)^{1/m}.$$

Theorem 3.13 will allow us to establish a useful result regarding the eigenvalues of any matrix A . Let B_1, B_2, \dots be a sequence of $m \times m$ matrices such that $B_n \rightarrow A$, as $n \rightarrow \infty$, and let $\delta(A, B_n)$ be as defined in Theorem 3.13. It follows from the fact that $B_n \rightarrow A$, as $n \rightarrow \infty$, that $\delta(A, B_n) \rightarrow 0$, as $n \rightarrow \infty$. Hence, if $\gamma_{1,n}, \dots, \gamma_{m,n}$ are the eigenvalues of B_n , then Theorem 3.13 tells us that

$$\max_{1 \leq i \leq m} \min_{1 \leq j \leq m} |\lambda_i - \gamma_{j,n}| \rightarrow 0,$$

as $n \rightarrow \infty$. In other words, if B_n is very close to A , then for each i , some j exists, such that $\gamma_{j,n}$ is close to λ_i , or more precisely, as $B_n \rightarrow A$, the eigenvalues of B_n are converging to those of A . This leads to Theorem 3.14.

Theorem 3.14 Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of the $m \times m$ matrix A . Then, for each i , λ_i is a continuous function of the elements of A .

Theorem 3.15 addresses the continuity of the eigenprojection $P_A(\lambda)$ of a symmetric matrix A . A detailed treatment of this problem, as well as the more general problem of the continuity of the eigenprojections of nonsymmetric matrices, can be found in Kato (1982).

Theorem 3.15 Suppose that A is an $m \times m$ symmetric matrix and λ is one of its eigenvalues. Then $P_A(\lambda)$, the eigenprojection associated with the eigenvalue λ , is a continuous function of the elements of A .

Example 3.11 Consider the matrix

$$A = \begin{bmatrix} 2 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

which clearly has the simple eigenvalue 2 and the repeated eigenvalue 1. Suppose that B_1, B_2, \dots is a sequence of 3×3 matrices such that $B_n \rightarrow A$, as $n \rightarrow \infty$.

Let $\gamma_{1,n} \geq \gamma_{2,n} \geq \gamma_{3,n}$ be the eigenvalues of B_n , whereas $\mathbf{x}_{1,n}$, $\mathbf{x}_{2,n}$, and $\mathbf{x}_{3,n}$ is a corresponding set of orthonormal eigenvectors. Theorem 3.14 implies that, as $n \rightarrow \infty$,

$$\gamma_{1,n} \rightarrow 2 \quad \text{and} \quad \gamma_{i,n} \rightarrow 1, \quad \text{for } i = 2, 3.$$

On the other hand, Theorem 3.15 implies that, as $n \rightarrow \infty$,

$$P_{1,n} \rightarrow P_A(2), \quad P_{2,n} \rightarrow P_A(1),$$

where

$$P_{1,n} = \mathbf{x}_{1,n} \mathbf{x}'_{1,n}, \quad P_{2,n} = \mathbf{x}_{2,n} \mathbf{x}'_{2,n} + \mathbf{x}_{3,n} \mathbf{x}'_{3,n}.$$

For instance, suppose that

$$B_n = \begin{bmatrix} 2 & 0 & n^{-1} \\ 0 & 1 & 0 \\ n^{-1} & 0 & 1 \end{bmatrix},$$

so that, clearly, $B_n \rightarrow A$. The characteristic equation of B_n simplifies to

$$\lambda^3 - 4\lambda^2 + (5 - n^{-2})\lambda - 2 + n^{-2} = (\lambda - 1)(\lambda^2 - 3\lambda + 2 - n^{-2}) = 0,$$

so that the eigenvalues of B_n are

$$1, \quad \frac{3}{2} - \frac{\sqrt{1 + 4n^{-2}}}{2}, \quad \frac{3}{2} + \frac{\sqrt{1 + 4n^{-2}}}{2},$$

which do converge to 1, 1, and 2, respectively. It is left as an exercise for the reader to verify that

$$P_{1,n} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = P_A(2), \quad P_{2,n} \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = P_A(1).$$

3.6 EXTREMAL PROPERTIES OF EIGENVALUES

One of the reasons that eigenvalues play a prominent role in many applications is because they can be expressed as maximum or minimum values of certain functions involving a quadratic form. In this section, we derive some of these extremal properties of eigenvalues.

Let A be a fixed $m \times m$ symmetric matrix, and consider the quadratic form $\mathbf{x}'A\mathbf{x}$ as a function of $\mathbf{x} \neq \mathbf{0}$. If α is a nonzero scalar, then $(\alpha\mathbf{x})'A(\alpha\mathbf{x}) = \alpha^2\mathbf{x}'A\mathbf{x}$, so that the quadratic form can be made arbitrarily small or large, depending on whether $\mathbf{x}'A\mathbf{x}$ is negative or positive, through the proper choice of α . Thus, any meaningful

study of the variational properties of $\mathbf{x}'A\mathbf{x}$ as we change \mathbf{x} will require the removal of the effect of scale changes in \mathbf{x} , which can be accomplished through the construction of what is commonly called the Rayleigh quotient given by

$$R(\mathbf{x}, A) = \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}}.$$

Note that $R(\alpha\mathbf{x}, A) = R(\mathbf{x}, A)$. Our first result involves the global maximization and minimization of $R(\mathbf{x}, A)$.

Theorem 3.16 Let A be a symmetric $m \times m$ matrix with ordered eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$. For any $m \times 1$ vector $\mathbf{x} \neq \mathbf{0}$,

$$\lambda_m \leq \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}} \leq \lambda_1, \quad (3.6)$$

and, in particular,

$$\lambda_m = \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}}, \quad \lambda_1 = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}}. \quad (3.7)$$

Proof. Let $A = X\Lambda X'$ be the spectral decomposition of A , where the columns of $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ are normalized eigenvectors of A and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$. Then, if $\mathbf{y} = X'\mathbf{x}$, we have

$$\frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}} = \frac{\mathbf{x}'X\Lambda X'\mathbf{x}}{\mathbf{x}X X'\mathbf{x}} = \frac{\mathbf{y}'\Lambda\mathbf{y}}{\mathbf{y}'\mathbf{y}} = \frac{\sum_{i=1}^m \lambda_i y_i^2}{\sum_{i=1}^m y_i^2},$$

so that (3.6) follows from the fact that

$$\lambda_m \sum_{i=1}^m y_i^2 \leq \sum_{i=1}^m \lambda_i y_i^2 \leq \lambda_1 \sum_{i=1}^m y_i^2.$$

Now (3.7) is verified by choices of \mathbf{x} for which the bounds in (3.6) are attained; for instance, the lower bound is attained with $\mathbf{x} = \mathbf{x}_m$, whereas the upper bound holds with $\mathbf{x} = \mathbf{x}_1$. \square

Note that, because for any nonnull \mathbf{x} , $\mathbf{z} = (\mathbf{x}'\mathbf{x})^{-1/2}\mathbf{x}$ is a unit vector, the minimization and maximization of $\mathbf{z}'A\mathbf{z}$ over all unit vectors \mathbf{z} will also yield λ_m and λ_1 , respectively; that is,

$$\lambda_m = \min_{\mathbf{z}'\mathbf{z}=1} \mathbf{z}'A\mathbf{z}, \quad \lambda_1 = \max_{\mathbf{z}'\mathbf{z}=1} \mathbf{z}'A\mathbf{z}.$$

Theorem 3.17 shows that each eigenvalue of a symmetric matrix A can be expressed as a constrained maximum or minimum of the Rayleigh quotient, $R(\mathbf{x}, A)$.

Theorem 3.17 Let A be an $m \times m$ symmetric matrix having eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ with $\mathbf{x}_1, \dots, \mathbf{x}_m$ being a corresponding set of orthonormal eigenvectors. For $h = 1, \dots, m$, define S_h and T_h to be the vector spaces spanned by the columns of $X_h = (\mathbf{x}_1, \dots, \mathbf{x}_h)$ and $Y_h = (\mathbf{x}_h, \dots, \mathbf{x}_m)$, respectively. Then

$$\lambda_h = \min_{\substack{\mathbf{x} \in S_h \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \min_{\substack{Y_{h+1}' \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}}$$

and

$$\lambda_h = \max_{\substack{\mathbf{x} \in T_h \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \max_{\substack{X_{h-1}' \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}}.$$

Proof. We will prove the result concerning the minimum; the proof for the maximum is similar. Let $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ and $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$. Note that, because $X'AX = \Lambda$ and $X'X = I_m$, it follows that $X_h'X_h = I_h$ and $X_h'AX_h = \Lambda_h$, where $\Lambda_h = \text{diag}(\lambda_1, \dots, \lambda_h)$. Now $\mathbf{x} \in S_h$ if and only if an $h \times 1$ vector \mathbf{y} exists, such that $\mathbf{x} = X_h\mathbf{y}$. Consequently,

$$\min_{\substack{\mathbf{x} \in S_h \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \min_{\mathbf{y} \neq \mathbf{0}} \frac{\mathbf{y}' X_h' A X_h \mathbf{y}}{\mathbf{y}' X_h' X_h \mathbf{y}} = \min_{\mathbf{y} \neq \mathbf{0}} \frac{\mathbf{y}' \Lambda_h \mathbf{y}}{\mathbf{y}' \mathbf{y}} = \lambda_h,$$

where the last equality follows from Theorem 3.16. The second version of the minimization follows immediately from the first and the fact that the null space of Y_{h+1}' is S_h . \square

Example 3.12 and Example 3.13 give some indication of how the extremal properties of eigenvalues make them important features in many applications.

Example 3.12 Suppose that the same m variables are measured on individuals from k different groups with the goal being to identify differences in the means for the k groups. Let the $m \times 1$ vectors $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ represent the k group mean vectors, and let $\boldsymbol{\mu} = (\boldsymbol{\mu}_1 + \dots + \boldsymbol{\mu}_k)/k$ be the average of these mean vectors. To investigate the differences in group means, we will use the deviations $(\boldsymbol{\mu}_i - \boldsymbol{\mu})$ from the average mean; in particular, we form the sum of squares and cross products matrix given by

$$A = \sum_{i=1}^k (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})'.$$

Note that for a particular unit vector \mathbf{x} , $\mathbf{x}' A \mathbf{x}$ will give a measure of the differences among the k groups in the direction \mathbf{x} ; a value of zero indicates the groups have

identical means in this direction, whereas increasingly large values of $\mathbf{x}'A\mathbf{x}$ indicate increasingly widespread differences in this same direction. If $\mathbf{x}_1, \dots, \mathbf{x}_m$ are normalized eigenvectors of A corresponding to its ordered eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$, then it follows from Theorem 3.16 and Theorem 3.17 that the greatest difference among the k groups, in terms of deviations from the overall mean, occurs in the direction given by \mathbf{x}_1 . Of all directions orthogonal to \mathbf{x}_1 , \mathbf{x}_2 gives the direction with the greatest difference among the k groups, and so on. If some of the eigenvalues are very small relative to the rest, then we will be able to effectively reduce the dimension of the problem. For example, suppose that $\lambda_3, \dots, \lambda_m$ are all very small relative to λ_1 and λ_2 . Then all substantial differences among the group means will be observed in the plane spanned by \mathbf{x}_1 and \mathbf{x}_2 . In Example 4.12, we will discuss the statistical procedure, called canonical variate analysis, which puts this sort of dimension reducing process into practice.

Example 3.13 In Example 3.12, the focus was on means. In this example, we will look at a procedure that concentrates on variances. This technique, called principal components analysis, was developed by Hotelling (1933). Some good references on this subject are Jackson (1991) and Jolliffe (2002). Let \mathbf{x} be an $m \times 1$ random vector having the covariance matrix Ω . Suppose that we wish to find the $m \times 1$ vector \mathbf{a}_1 so as to make the variance of $\mathbf{a}_1'\mathbf{x}$ as large as possible. However, from Section 1.13, we know that

$$\text{var}(\mathbf{a}_1'\mathbf{x}) = \mathbf{a}_1'\{\text{var}(\mathbf{x})\}\mathbf{a}_1 = \mathbf{a}_1'\Omega\mathbf{a}_1. \quad (3.8)$$

Clearly, we can make this variance arbitrarily large by taking $\mathbf{a}_1 = \alpha\mathbf{c}$ for some scalar α and some vector $\mathbf{c} \neq \mathbf{0}$, and then let $\alpha \rightarrow \infty$. We will remove this effect of the scale of \mathbf{a}_1 by imposing a constraint. For example, we may consider maximizing (3.8) over all choices of \mathbf{a}_1 satisfying $\mathbf{a}_1'\mathbf{a}_1 = 1$. In this case, we are searching for the one direction in R^m , that is, the line, for which the variability of observations of \mathbf{x} projected onto that line is maximized. It follows from Theorem 3.16 that this direction is given by the normalized eigenvector of Ω corresponding to its largest eigenvalue. Suppose we also wish to find a second direction, given by \mathbf{a}_2 and orthogonal to \mathbf{a}_1 , where $\mathbf{a}_2'\mathbf{a}_2 = 1$ and $\text{var}(\mathbf{a}_2'\mathbf{x})$ is maximized. From Theorem 3.17, this second direction is given by the normalized eigenvector of Ω corresponding to its second largest eigenvalue. Continuing in this fashion, we would obtain m directions identified by the set $\mathbf{a}_1, \dots, \mathbf{a}_m$ of orthonormal eigenvectors of Ω . Effectively, what we will have done is to find a rotation of the original axes to a new set of orthogonal axes, where each successive axis is selected so as to maximize the dispersion among the \mathbf{x} observations along that axis. Note that the components of the transformed vector $(\mathbf{a}_1'\mathbf{x}, \dots, \mathbf{a}_m'\mathbf{x})'$, which are called the principal components of Ω , are uncorrelated because for $i \neq j$,

$$\text{cov}(\mathbf{a}_i'\mathbf{x}, \mathbf{a}_j'\mathbf{x}) = \mathbf{a}_i'\Omega\mathbf{a}_j = \mathbf{a}_i'(\lambda_j\mathbf{a}_j) = \lambda_j\mathbf{a}_i'\mathbf{a}_j = 0.$$

For some specific examples, first consider the 4×4 covariance matrix given by

$$\Omega = \begin{bmatrix} 4.65 & 4.35 & 0.55 & 0.45 \\ 4.35 & 4.65 & 0.45 & 0.55 \\ 0.55 & 0.45 & 4.65 & 4.35 \\ 0.45 & 0.55 & 4.35 & 4.65 \end{bmatrix}.$$

The eigenvalues of Ω are 10, 8, 0.4, and 0.2, so the first two eigenvalues account for a large proportion, actually $18/18.6 = 0.97$, of the total variability of \mathbf{x} , which means that although observations of \mathbf{x} would appear as points in R^4 , almost all of the dispersion among these points will be confined to a plane. This plane is spanned by the first two normalized eigenvectors of Ω , $(0.5, 0.5, 0.5, 0.5)'$ and $(0.5, 0.5, -0.5, -0.5)'$. As a second illustration, consider a covariance matrix such as

$$\Omega = \begin{bmatrix} 59 & 5 & 2 \\ 5 & 35 & -10 \\ 2 & -10 & 56 \end{bmatrix},$$

which has a repeated eigenvalue; specifically the eigenvalues are 60 and 30 with multiplicities 2 and 1, respectively. Since the largest eigenvalue of Ω is repeated, there is no one direction \mathbf{a}_1 that maximizes $\text{var}(\mathbf{a}_1' \mathbf{x})$. Instead, the dispersion of \mathbf{x} observations is the same in all directions in the plane given by the eigenspace $S_\Omega(60)$, which is spanned by the vectors $(1, 1, -2)'$ and $(2, 0, 1)'$. Consequently, a scatter plot of \mathbf{x} observations would produce a circular pattern of points in this plane.

Theorem 3.18, known as the Courant–Fischer min–max theorem, gives alternative expressions for the intermediate eigenvalues of A as constrained minima and maxima of the Rayleigh quotient $R(\mathbf{x}, A)$.

Theorem 3.18 Let A be an $m \times m$ symmetric matrix having eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. For $h = 1, \dots, m$, let B_h be any $m \times (h-1)$ matrix and C_h any $m \times (m-h)$ matrix satisfying $B_h' B_h = I_{h-1}$ and $C_h' C_h = I_{m-h}$. Then

$$\lambda_h = \min_{B_h} \max_{\substack{B_h' \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \quad (3.9)$$

as well as

$$\lambda_h = \min_{C_h} \max_{\substack{C_h' \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \quad (3.10)$$

Proof. We first prove the min–max result given by (3.9). Let $X_h = (\mathbf{x}_1, \dots, \mathbf{x}_h)$, where $\mathbf{x}_1, \dots, \mathbf{x}_h$ is a set of orthonormal eigenvectors of A , corresponding to

the eigenvalues $\lambda_1, \dots, \lambda_h$. Since X_{h-1} is an $m \times (h-1)$ matrix satisfying $X'_{h-1}X_{h-1} = I_{h-1}$, it follows that

$$\min_{B_h} \max_{\substack{B'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \leq \max_{\substack{X'_{h-1} \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \lambda_h, \quad (3.11)$$

where the equality follows from Theorem 3.17. Now for arbitrary B_h satisfying $B'_h B_h = I_{h-1}$, the matrix $B'_h X_h$ is $(h-1) \times h$, so that the columns must be linearly dependent. Consequently, we can find an $h \times 1$ nonnull vector \mathbf{y} such that $B'_h X_h \mathbf{y} = \mathbf{0}$. Since $X_h \mathbf{y}$ is one choice for \mathbf{x} , we find that

$$\max_{\substack{B'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \geq \frac{\mathbf{y}' X'_h A X_h \mathbf{y}}{\mathbf{y}' X'_h X_h \mathbf{y}} = \frac{\mathbf{y}' \Lambda_h \mathbf{y}}{\mathbf{y}' \mathbf{y}} \geq \lambda_h, \quad (3.12)$$

where $\Lambda_h = \text{diag}(\lambda_1, \dots, \lambda_h)$ and the last inequality follows from (3.6). Minimizing (3.12) over all choices of B_h gives

$$\min_{B_h} \max_{\substack{B'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \geq \lambda_h.$$

This, along with (3.11), proves (3.9). The proof of (3.10) is along the same lines. Let $Y_h = (\mathbf{x}_h, \dots, \mathbf{x}_m)$, where $\mathbf{x}_h, \dots, \mathbf{x}_m$ is a set of orthonormal eigenvectors of A , corresponding to the eigenvalues $\lambda_h, \dots, \lambda_m$. Since Y_{h+1} is an $m \times (m-h)$ matrix satisfying $Y'_{h+1} Y_{h+1} = I_{m-h}$, it follows that

$$\max_{C_h} \min_{\substack{C'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \geq \min_{\substack{Y'_{h+1} \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} = \lambda_h, \quad (3.13)$$

where the equality follows from Theorem 3.17. For an arbitrary C_h satisfying $C'_h C_h = I_{m-h}$, the matrix $C'_h Y_h$ is $(m-h) \times (m-h+1)$, so the columns of $C'_h Y_h$ must be linearly dependent. Thus, an $(m-h+1) \times 1$ nonnull vector \mathbf{y} exists, satisfying $C'_h Y_h \mathbf{y} = \mathbf{0}$. Since $Y_h \mathbf{y}$ is one choice for \mathbf{x} , we have

$$\min_{\substack{C'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \leq \frac{\mathbf{y}' Y'_h A Y_h \mathbf{y}}{\mathbf{y}' Y'_h Y_h \mathbf{y}} = \frac{\mathbf{y}' \Delta_h \mathbf{y}}{\mathbf{y}' \mathbf{y}} \leq \lambda_h, \quad (3.14)$$

where $\Delta_h = \text{diag}(\lambda_h, \dots, \lambda_m)$ and the last inequality follows from (3.6). Maximizing (3.14) over all choices of C_h yields

$$\max_{C_h} \min_{\substack{C'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \leq \lambda_h,$$

which together with (3.13) establishes (3.10). \square

Corollary 3.18.1 Let A be an $m \times m$ symmetric matrix having eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$. For $h = 1, \dots, m$, let B_h be any $m \times (h-1)$ matrix and C_h be any $m \times (m-h)$ matrix. Then

$$\lambda_h \leq \max_{\substack{B'_h \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}},$$

and

$$\lambda_h \leq \min_{\substack{C'_h \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}}.$$

Proof. If $B'_h B_h = I_{h-1}$ and $C'_h C_h = I_{m-h}$, then the two inequalities follow directly from Theorem 3.18. We need to establish them for arbitrary B_h and C_h . When $B'_h B_h = I_{h-1}$, the set $S_{B_h} = \{\mathbf{x} : \mathbf{x} \in R^m, B'_h \mathbf{x} = 0\}$ is the orthogonal complement of the vector space that has the columns of B_h as an orthonormal basis. Thus, the first inequality holds when maximizing over all $\mathbf{x} \neq 0$ in any $(m-h+1)$ -dimensional vector subspace of R^m . Consequently, this inequality also will hold for any $m \times (h-1)$ matrix B_h because, in this case, $\text{rank}(B_h) \leq h-1$ guarantees that the maximization is over a vector subspace of dimension at least $m-h+1$. A similar argument applies to the second inequality. \square

The proof of the following extension of Theorem 3.18 is left to the reader as an exercise.

Corollary 3.18.2 Let A be an $m \times m$ symmetric matrix having eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$, and let i_1, \dots, i_k be integers satisfying $1 \leq i_1 < \dots < i_k \leq m$. Define the matrices B_{i_1}, \dots, B_{i_k} , such that B_{i_j} is $m \times (i_j-1)$, $B'_{i_j} B_{i_j} = I_{i_j-1}$, and $B_{i_h} B'_{i_h} B_{i_j} = B_{i_j}$ for $h = j+1, \dots, k$. Define the matrices C_{i_1}, \dots, C_{i_k} , such that C_{i_j} is $m \times (m-i_j)$, $C'_{i_j} C_{i_j} = I_{m-i_j}$, and $C_{i_h} C'_{i_h} C_{i_j} = C_{i_j}$ for $h = 1, \dots, j$. Then

$$\sum_{j=1}^k \lambda_{i_j} = \min_{B_{i_1}, \dots, B_{i_k}} \max_{\substack{B'_{i_1} \mathbf{x}_1 = \dots = B'_{i_k} \mathbf{x}_k = 0 \\ \mathbf{x}_1 \neq 0, \dots, \mathbf{x}_k \neq 0 \\ \mathbf{x}'_h \mathbf{x}_l = 0, h \neq l}} \sum_{j=1}^k \frac{\mathbf{x}'_j A \mathbf{x}_j}{\mathbf{x}'_j \mathbf{x}_j}$$

and

$$\sum_{j=1}^k \lambda_{i_j} = \max_{C_{i_1}, \dots, C_{i_k}} \max_{\substack{C'_{i_1} \mathbf{x}_1 = \dots = C'_{i_k} \mathbf{x}_k = 0 \\ \mathbf{x}_1 \neq 0, \dots, \mathbf{x}_k \neq 0 \\ \mathbf{x}'_h \mathbf{x}_l = 0, h \neq l}} \sum_{j=1}^k \frac{\mathbf{x}'_j A \mathbf{x}_j}{\mathbf{x}'_j \mathbf{x}_j}.$$

3.7 ADDITIONAL RESULTS CONCERNING EIGENVALUES OF SYMMETRIC MATRICES

Let A be an $m \times m$ symmetric matrix and H be an $m \times h$ matrix satisfying $H'H = I_h$. In some situations, it is of interest to compare the eigenvalues of A with those of $H'AH$. Some comparisons follow immediately from Theorem 3.18. For instance, it is easily verified that from (3.9), we have

$$\lambda_1(H'AH) \geq \lambda_{m-h+1}(A),$$

and from (3.10) we have

$$\lambda_h(H'AH) \leq \lambda_h(A).$$

Theorem 3.19, known as the Poincaré separation theorem (Poincaré, 1890; see also Fan, 1949), provides some inequalities involving the eigenvalues of A and $H'AH$ in addition to the two given above.

Theorem 3.19 Let A be an $m \times m$ symmetric matrix and H be an $m \times h$ matrix satisfying $H'H = I_h$. Then, for $i = 1, \dots, h$, it follows that

$$\lambda_{m-h+i}(A) \leq \lambda_i(H'AH) \leq \lambda_i(A).$$

Proof. To establish the lower bound on $\lambda_i(H'AH)$, let $Y_n = (\mathbf{x}_n, \dots, \mathbf{x}_m)$, where $n = m - h + i + 1$, and $\mathbf{x}_1, \dots, \mathbf{x}_m$ is a set of orthonormal eigenvectors of A corresponding to the eigenvalues $\lambda_1(A) \geq \dots \geq \lambda_m(A)$. Then it follows that

$$\begin{aligned} \lambda_{m-h+i}(A) &= \lambda_{n-1}(A) = \min_{\substack{Y'_n \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \leq \min_{\substack{Y'_n \mathbf{x} = 0 \\ \mathbf{x} = H \mathbf{y} \\ \mathbf{y} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \\ &= \min_{\substack{Y'_n H \mathbf{y} = 0 \\ \mathbf{y} \neq 0}} \frac{\mathbf{y}' H' A H \mathbf{y}}{\mathbf{y}' \mathbf{y}} \leq \lambda_{h-(m-n+1)}(H'AH) \\ &= \lambda_i(H'AH), \end{aligned}$$

where the second equality follows from Theorem 3.17. The last inequality follows from Corollary 3.18.1, after noting that the order of $H'AH$ is h and $Y'_n H$ is $(m - n + 1) \times h$. To prove the upper bound for $\lambda_i(H'AH)$, let $X_{i-1} = (\mathbf{x}_1, \dots, \mathbf{x}_{i-1})$, and note that

$$\begin{aligned} \lambda_i(A) &= \max_{\substack{X'_{i-1} \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \geq \max_{\substack{X'_{i-1} \mathbf{x} = 0 \\ \mathbf{x} = H \mathbf{y} \\ \mathbf{y} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \\ &= \max_{\substack{X'_{i-1} H \mathbf{y} = 0 \\ \mathbf{y} \neq 0}} \frac{\mathbf{y}' H' A H \mathbf{y}}{\mathbf{y}' \mathbf{y}} \geq \lambda_i(H'AH), \end{aligned}$$

where the first equality follows from Theorem 3.17 and the final inequality follows from Corollary 3.18.1. \square

Theorem 3.19 can be used to prove Theorem 3.20.

Theorem 3.20 Let A be an $m \times m$ symmetric matrix, and let A_k be its leading $k \times k$ principal submatrix; that is, A_k is the matrix obtained by deleting the last $m - k$ rows and columns of A . Then, for $i = 1, \dots, k$,

$$\lambda_{m-i+1}(A) \leq \lambda_{k-i+1}(A_k) \leq \lambda_{k-i+1}(A).$$

Theorem 3.21, sometimes referred to as Weyl's Theorem, gives inequalities relating the eigenvalues of two symmetric matrices to those of the sum of the matrices.

Theorem 3.21 Let A and B be $m \times m$ symmetric matrices. Then for $h = 1, \dots, m$,

$$\lambda_h(A) + \lambda_m(B) \leq \lambda_h(A + B) \leq \lambda_h(A) + \lambda_1(B).$$

Proof. Let B_h be an $m \times (h - 1)$ matrix satisfying $B_h' B_h = I_{h-1}$. Then using (3.9), we have

$$\begin{aligned} \lambda_h(A + B) &= \min_{B_h} \max_{\substack{B_h' x = 0 \\ x \neq 0}} \frac{x'(A + B)x}{x'x} \\ &= \min_{B_h} \max_{\substack{B_h' x = 0 \\ x \neq 0}} \left(\frac{x'Ax}{x'x} + \frac{x'Bx}{x'x} \right) \\ &\geq \min_{B_h} \max_{\substack{B_h' x = 0 \\ x \neq 0}} \left(\frac{x'Ax}{x'x} + \lambda_m(B) \right) \\ &= \min_{B_h} \max_{\substack{B_h' x = 0 \\ x \neq 0}} \left(\frac{x'Ax}{x'x} \right) + \lambda_m(B) \\ &= \lambda_h(A) + \lambda_m(B), \end{aligned}$$

where the inequality was introduced from an application of (3.6), whereas the final equality used (3.9). The upper bound is obtained in a similar manner by using (3.10). \square

The inequalities given in Theorem 3.21 can be generalized. Before obtaining these generalized inequalities, we first give some inequalities relating the eigenvalues of $A + B$ to those of A when we have some information about the rank of B .

Theorem 3.22 Let A and B be $m \times m$ symmetric matrices, and suppose that $\text{rank}(B) \leq r$. Then for $h = 1, \dots, m - r$,

- (a) $\lambda_{h+r}(A) \leq \lambda_h(A + B)$,
- (b) $\lambda_{h+r}(A + B) \leq \lambda_h(A)$.

Proof. Since B is symmetric and has rank at most r , it can be expressed as

$$B = \sum_{i=1}^r \gamma_i \mathbf{y}_i \mathbf{y}_i',$$

where $\mathbf{y}_1, \dots, \mathbf{y}_r$ are orthonormal vectors. Let B_h and B_{h+r} be $m \times (h-1)$ and $m \times (h+r-1)$ matrices satisfying $B_h' B_h = I_{h-1}$ and $B_{h+r}' B_{h+r} = I_{h+r-1}$, and define $Y_r = (\mathbf{y}_1, \dots, \mathbf{y}_r)$ and $B_* = (B_h, Y_r)$. Using (3.9), if $h = 1, \dots, m - r$, we have

$$\begin{aligned} \lambda_h(A + B) &= \min_{B_h} \max_{\substack{B_h' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'(A + B)\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\ &\geq \min_{B_h} \max_{\substack{B_h' \mathbf{x} = 0, Y_r' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'(A + B)\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\ &= \min_{B_h} \max_{\substack{B_h' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \\ &= \min_{B_h Y_r = (0)} \max_{\substack{B_h' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \\ &\geq \min_{B_{h+r}} \max_{\substack{B_{h+r}' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \\ &= \lambda_{h+r}(A) \end{aligned}$$

The third equality can be justified as follows. First, the minimum can be restricted to B_h for which B_* is full rank, because if B_* is not full rank, then $\{\mathbf{x} : B_*' \mathbf{x} = \mathbf{0}\}$ will contain subspaces of the form $\{\mathbf{x} : B_h' \mathbf{x} = \mathbf{0}\}$ for choices of B_h for which B_* is full rank. Secondly, for choices of B_h for which B_* is full rank, the null space of B_*' , $\{\mathbf{x} : B_*' \mathbf{x} = \mathbf{0}\}$, will be identical to the null space of B_h' for a choice of B_h that has $B_h' Y_r = (0)$. This establishes (a). We obtain (b) in a similar fashion by using (3.10). Let C_h and C_{h+r} be $m \times (m-h)$ and $m \times (m-h-r)$ matrices satisfying $C_h' C_h = I_{m-h}$ and $C_{h+r}' C_{h+r} = I_{m-h-r}$, and define $C_* = (C_{h+r}, Y_r)$. Then for $h = 1, \dots, m - r$,

$$\lambda_{h+r}(A + B) = \max_{C_{h+r}} \min_{\substack{C_{h+r}' \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'(A + B)\mathbf{x}}{\mathbf{x}'\mathbf{x}}$$

$$\begin{aligned}
&\leq \max_{C_{h+r}} \min_{\substack{C'_{h+r} \\ \mathbf{x}=0, \mathbf{y}'_r \mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'(A+B)\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\
&= \max_{C_{h+r}} \min_{\substack{C'_h \mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\
&= \max_{C'_{h+r} \mathbf{y}_r=(0)} \min_{\substack{C'_h \mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\
&\leq \max_{C_h} \min_{\substack{C'_h \mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'\mathbf{x}} \\
&= \lambda_h(A)
\end{aligned}$$

and so the proof is complete. \square

We are now ready to give a generalization of the inequalities given in Theorem 3.21.

Theorem 3.23 Let A and B be $m \times m$ symmetric matrices, and let h and i be integers between 1 and m inclusive. Then

- (a) $\lambda_{h+i-1}(A+B) \leq \lambda_h(A) + \lambda_i(B)$, if $h+i \leq m+1$,
- (b) $\lambda_{h+i-m}(A+B) \geq \lambda_h(A) + \lambda_i(B)$, if $h+i \geq m+1$.

Proof. Let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be a set of orthonormal eigenvectors corresponding to the eigenvalues $\lambda_1(A) \geq \dots \geq \lambda_m(A)$ and $\mathbf{y}_1, \dots, \mathbf{y}_m$ be a set of orthonormal eigenvectors corresponding to the eigenvalues $\lambda_1(B) \geq \dots \geq \lambda_m(B)$. If we define $A_h = \sum_{j=1}^{h-1} \lambda_j(A) \mathbf{x}_j \mathbf{x}'_j$ and $B_i = \sum_{j=1}^{i-1} \lambda_j(B) \mathbf{y}_j \mathbf{y}'_j$, then clearly $\text{rank}(A_h) \leq h-1$, $\text{rank}(B_i) \leq i-1$, and $\text{rank}(A_h + B_i) \leq h+i-2$. Thus, applying Theorem 3.22(a) with $h=1$ and $r=h+i-2$, we have

$$\begin{aligned}
\lambda_1(A - A_h + B - B_i) &= \lambda_1((A+B) - (A_h + B_i)) \\
&\geq \lambda_{1+h+i-2}(A+B) \\
&= \lambda_{h+i-1}(A+B).
\end{aligned} \tag{3.15}$$

Also

$$\lambda_1(A - A_h + B - B_i) \leq \lambda_1(A - A_h) + \lambda_1(B - B_i) \tag{3.16}$$

from Theorem 3.21. Note that

$$\lambda_1(A - A_h) = \lambda_h(A), \quad \lambda_1(B - B_i) = \lambda_i(B) \tag{3.17}$$

because

$$A - A_h = \sum_{j=h}^m \lambda_j(A) \mathbf{x}_j \mathbf{x}_j', \quad B - B_i = \sum_{j=i}^m \lambda_j(B) \mathbf{y}_j \mathbf{y}_j'.$$

Combining (3.15), (3.16), and (3.17), we get

$$\begin{aligned} \lambda_h(A) + \lambda_i(B) &= \lambda_1(A - A_h) + \lambda_1(B - B_i) \\ &\geq \lambda_1(A - A_h + B - B_i) \\ &= \lambda_1((A + B) - (A_h + B_i)) \\ &\geq \lambda_{h+i-1}(A + B), \end{aligned}$$

so that we have proven the inequality given in (a). The inequality in (b) can be obtained by applying the inequality in (a) to $-A$ and $-B$; that is,

$$\lambda_{h+i-1}(-A - B) \leq \lambda_h(-A) + \lambda_i(-B),$$

which can be re-expressed as

$$-\lambda_{m-(h+i-1)+1}(A + B) \leq -\lambda_{m-h+1}(A) - \lambda_{m-i+1}(B),$$

or equivalently,

$$\lambda_{m-h-i+2}(A + B) \geq \lambda_{m-h+1}(A) + \lambda_{m-i+1}(B).$$

The last inequality is identical to the one given in (b) because if we let $k = m - h + 1$ and $l = m - i + 1$, then $k + l - m = m - h - i + 2$. \square

The inequalities given in the preceding theorems can be used to obtain bounds on sums of eigenvalues. For instance, from Theorem 3.21 it immediately follows that

$$\sum_{h=1}^k \lambda_h(A) + k\lambda_m(B) \leq \sum_{h=1}^k \lambda_h(A + B) \leq \sum_{h=1}^k \lambda_h(A) + k\lambda_1(B).$$

Theorem 3.24, which is due to Wielandt (1955), provides tighter bounds on the sums of eigenvalues of $A + B$.

Theorem 3.24 Let A and B be $m \times m$ symmetric matrices, and let i_1, \dots, i_k be integers satisfying $1 \leq i_1 < \dots < i_k \leq m$. Then for $k = 1, \dots, m$,

$$\sum_{j=1}^k \{\lambda_{i_j}(A) + \lambda_{m-k+j}(B)\} \leq \sum_{j=1}^k \lambda_{i_j}(A + B) \leq \sum_{j=1}^k \{\lambda_{i_j}(A) + \lambda_j(B)\}.$$

Proof. Note that it follows from Corollary 3.18.2 that there are particular matrices C_{i_1}, \dots, C_{i_k} , such that C_{i_j} is $m \times (m - i_j)$, $C'_{i_j} C_{i_j} = I_{m-i_j}$, $C'_{i_h} C'_{i_h} C_{i_j} = C_{i_j}$ for $h = 1, \dots, j$, and

$$\sum_{j=1}^k \lambda_{i_j}(A + B) = \min_{\substack{C'_{i_1} \mathbf{x}_1 = \dots = C'_{i_k} \mathbf{x}_k = \mathbf{0} \\ \mathbf{x}_1 \neq \mathbf{0}, \dots, \mathbf{x}_k \neq \mathbf{0} \\ \mathbf{x}'_h \mathbf{x}_l = 0, h \neq l}} \sum_{j=1}^k \frac{\mathbf{x}'_j (A + B) \mathbf{x}_j}{\mathbf{x}'_j \mathbf{x}_j}. \quad (3.18)$$

Let $\mathbf{y}_1, \dots, \mathbf{y}_k$ be $m \times 1$ unit vectors such that $\mathbf{y}'_h \mathbf{y}_l = 0$ for $h \neq l$, $C'_{i_1} \mathbf{y}_1 = \dots = C'_{i_k} \mathbf{y}_k = \mathbf{0}$, and

$$\sum_{j=1}^k \mathbf{y}'_j A \mathbf{y}_j = \min_{\substack{C'_{i_1} \mathbf{x}_1 = \dots = C'_{i_k} \mathbf{x}_k = \mathbf{0} \\ \mathbf{x}_1 \neq \mathbf{0}, \dots, \mathbf{x}_k \neq \mathbf{0} \\ \mathbf{x}'_h \mathbf{x}_l = 0, h \neq l}} \sum_{j=1}^k \frac{\mathbf{x}'_j A \mathbf{x}_j}{\mathbf{x}'_j \mathbf{x}_j}. \quad (3.19)$$

It follows from (3.18) that

$$\begin{aligned} \sum_{j=1}^k \lambda_{i_j}(A + B) &\leq \sum_{j=1}^k \mathbf{y}'_j (A + B) \mathbf{y}_j \\ &= \sum_{j=1}^k \mathbf{y}'_j A \mathbf{y}_j + \sum_{j=1}^k \mathbf{y}'_j B \mathbf{y}_j. \end{aligned} \quad (3.20)$$

Since the \mathbf{y}_j 's were chosen to satisfy (3.19), a direct application of Corollary 3.18.2 yields

$$\sum_{j=1}^k \mathbf{y}'_j A \mathbf{y}_j \leq \sum_{j=1}^k \lambda_{i_j}(A). \quad (3.21)$$

Let $\mathbf{y}_{k+1}, \dots, \mathbf{y}_m$ be unit vectors such that $\mathbf{y}_1, \dots, \mathbf{y}_m$ is an orthonormal set of vectors, and define the $m \times (m - i)$ matrix $C_{*i} = (\mathbf{y}_{i+1}, \dots, \mathbf{y}_m)$ for $i = 1, \dots, k$. Then it follows that

$$\sum_{j=1}^k \mathbf{y}'_j B \mathbf{y}_j = \min_{\substack{C'_{*1} \mathbf{x}_1 = \dots = C'_{*k} \mathbf{x}_k = \mathbf{0} \\ \mathbf{x}_1 \neq \mathbf{0}, \dots, \mathbf{x}_k \neq \mathbf{0} \\ \mathbf{x}'_h \mathbf{x}_l = 0, h \neq l}} \sum_{j=1}^k \frac{\mathbf{x}'_j B \mathbf{x}_j}{\mathbf{x}'_j \mathbf{x}_j},$$

and so another application of Corollary 3.18.2 leads to

$$\sum_{j=1}^k \mathbf{y}'_j B \mathbf{y}_j \leq \sum_{j=1}^k \lambda_j(B). \quad (3.22)$$

Using (3.21) and (3.22) in (3.20), we then get the required upper bound. The lower bound is established in a similar fashion by using the min–max identity given in Corollary 3.18.2. \square

Many applications utilizing Theorem 3.24 will involve the sum of the k largest eigenvalues of $A + B$. This special case of Theorem 3.24 is highlighted in the following corollary.

Corollary 3.24.1 Let A and B be $m \times m$ symmetric matrices. Then for $k = 1, \dots, m$,

$$\sum_{i=1}^k \lambda_i(A) + \sum_{i=1}^k \lambda_{m-k+i}(B) \leq \sum_{i=1}^k \lambda_i(A+B) \leq \sum_{i=1}^k \lambda_i(A) + \sum_{i=1}^k \lambda_i(B).$$

Some additional results regarding eigenvalues can be found in Bellman (1970) and Horn and Johnson (2013).

3.8 NONNEGATIVE DEFINITE MATRICES

In Chapter 1, the conditions for a symmetric matrix A to be a positive definite or positive semidefinite matrix were given in terms of the possible values of the quadratic form $\mathbf{x}'A\mathbf{x}$. We now show that these conditions also can be expressed in terms of the eigenvalues of A .

Theorem 3.25 Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of the $m \times m$ symmetric matrix A . Then

- (a) A is positive definite if and only if $\lambda_i > 0$ for all i ,
- (b) A is positive semidefinite if and only if $\lambda_i \geq 0$ for all i and $\lambda_i = 0$ for at least one i .

Proof. Let the columns of $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ be a set of orthonormal eigenvectors of A corresponding to the eigenvalues $\lambda_1, \dots, \lambda_m$, so that $A = X\Lambda X'$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$. If A is positive definite, then $\mathbf{x}'A\mathbf{x} > 0$ for all $\mathbf{x} \neq \mathbf{0}$, so in particular, choosing $\mathbf{x} = \mathbf{x}_i$, we have

$$\mathbf{x}_i' A \mathbf{x}_i = \mathbf{x}_i' (\lambda_i \mathbf{x}_i) = \lambda_i \mathbf{x}_i' \mathbf{x}_i = \lambda_i > 0.$$

Conversely, if $\lambda_i > 0$ for all i , then for any $\mathbf{x} \neq \mathbf{0}$ define $\mathbf{y} = X'\mathbf{x}$, and note that

$$\mathbf{x}' A \mathbf{x} = \mathbf{x}' X \Lambda X' \mathbf{x} = \mathbf{y}' \Lambda \mathbf{y} = \sum_{i=1}^m y_i^2 \lambda_i \quad (3.23)$$

has to be positive because the λ_i 's are positive and at least one of the y_i^2 's is positive because $\mathbf{y} \neq \mathbf{0}$. This proves (a). By a similar argument, we find that A is nonnegative definite if and only if $\lambda_i \geq 0$ for all i . Thus, to prove (b), we only need to prove that $\mathbf{x}'A\mathbf{x} = 0$ for some $\mathbf{x} \neq \mathbf{0}$ if and only if at least one $\lambda_i = 0$. It follows from (3.23) that if $\mathbf{x}'A\mathbf{x} = 0$, then $\lambda_i = 0$ for every i for which $y_i^2 > 0$. On the other hand, if for some i , $\lambda_i = 0$, then $\mathbf{x}'_i A \mathbf{x}_i = \lambda_i = 0$. \square

Since a square matrix is singular if and only if it has a zero eigenvalue, it follows immediately from Theorem 3.25 that positive definite matrices are nonsingular, whereas positive semidefinite matrices are singular.

Example 3.14 Consider the ordinary least squares estimator $\hat{\beta} = (X'X)^{-1}X'\mathbf{y}$ of β in the model

$$\mathbf{y} = X\beta + \epsilon,$$

where β is $(k+1) \times 1$, $E(\epsilon) = \mathbf{0}$ and $\text{var}(\epsilon) = \sigma^2 I_N$. For an arbitrary $(k+1) \times 1$ vector \mathbf{c} , we will prove that $\mathbf{c}'\hat{\beta}$ is the best linear unbiased estimator of $\mathbf{c}'\beta$; an estimator t is an unbiased estimator of $\mathbf{c}'\beta$ if $E(t) = \mathbf{c}'\beta$. Clearly, $\mathbf{c}'\hat{\beta}$ is unbiased because $E(\epsilon) = \mathbf{0}$ implies that

$$\begin{aligned} E(\mathbf{c}'\hat{\beta}) &= \mathbf{c}'(X'X)^{-1}X'E(\mathbf{y}) \\ &= \mathbf{c}'(X'X)^{-1}X'X\beta \\ &= \mathbf{c}'\beta. \end{aligned}$$

To show that $\mathbf{c}'\hat{\beta}$ is the best linear unbiased estimator, we must show that it has variance at least as small as the variance of any other linear unbiased estimator of $\mathbf{c}'\beta$. Let $\mathbf{a}'\mathbf{y}$ be an arbitrary linear unbiased estimator of $\mathbf{c}'\beta$, so that

$$\mathbf{c}'\beta = E(\mathbf{a}'\mathbf{y}) = \mathbf{a}'E(\mathbf{y}) = \mathbf{a}'X\beta,$$

regardless of the value of the vector β . However, this implies that

$$\mathbf{c}' = \mathbf{a}'X.$$

Now we saw in Example 3.7 that $\text{var}(\hat{\beta}) = \sigma^2(X'X)^{-1}$, so

$$\begin{aligned} \text{var}(\mathbf{c}'\hat{\beta}) &= \mathbf{c}'\{\text{var}(\hat{\beta})\}\mathbf{c} = \mathbf{c}'\{\sigma^2(X'X)^{-1}\}\mathbf{c} \\ &= \sigma^2\mathbf{a}'X(X'X)^{-1}X'\mathbf{a}, \end{aligned}$$

whereas

$$\text{var}(\mathbf{a}'\mathbf{y}) = \mathbf{a}'\{\text{var}(\mathbf{y})\}\mathbf{a} = \mathbf{a}'\{\sigma^2 I_N\}\mathbf{a} = \sigma^2\mathbf{a}'\mathbf{a}.$$

Thus, the difference in their variances is

$$\begin{aligned}\text{var}(\mathbf{a}'\mathbf{y}) - \text{var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) &= \sigma^2\mathbf{a}'\mathbf{a} - \sigma^2\mathbf{a}'X(X'X)^{-1}X'\mathbf{a} \\ &= \sigma^2\mathbf{a}'(I_N - X(X'X)^{-1}X')\mathbf{a}.\end{aligned}$$

However,

$$\{I_N - X(X'X)^{-1}X'\}^2 = \{I_N - X(X'X)^{-1}X'\},$$

and so using Theorem 3.4, we find that each of the eigenvalues of $I_N - X(X'X)^{-1}X'$ must be 0 or 1. Thus, from Theorem 3.25, we see that $I_N - X(X'X)^{-1}X'$ is nonnegative definite, and so

$$\text{var}(\mathbf{a}'\mathbf{y}) - \text{var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) \geq 0,$$

as is required.

Symmetric matrices are often obtained as the result of a transpose product; that is, if T is an $m \times n$ matrix, then both $T'T$ and TT' are symmetric matrices. The following two theorems show that their eigenvalues are nonnegative and their positive eigenvalues are equal.

Theorem 3.26 Let T be an $m \times n$ matrix with $\text{rank}(T) = r$. Then $T'T$ has r positive eigenvalues. It is positive definite if $r = n$ and positive semidefinite if $r < n$.

Proof. For any nonnull $n \times 1$ vector \mathbf{x} , let $\mathbf{y} = T\mathbf{x}$. Then clearly

$$\mathbf{x}'T'T\mathbf{x} = \mathbf{y}'\mathbf{y} = \sum_{i=1}^m y_i^2$$

is nonnegative, so $T'T$ is nonnegative definite, and thus, by Theorem 3.25, all of its eigenvalues are nonnegative. If \mathbf{x} is an eigenvector of $T'T$ corresponding to a zero eigenvalue, then the equation above must equal zero, and this can only happen if $\mathbf{y} = T\mathbf{x} = \mathbf{0}$. Since $\text{rank}(T) = r$, we can find a set of $n - r$ linearly independent \mathbf{x} 's satisfying $T\mathbf{x} = \mathbf{0}$, that is, any basis of the null space of T , and so the number of zero eigenvalues of $T'T$ is equal to $n - r$. The result now follows. \square

Theorem 3.27 Let T be an $m \times n$ matrix, with $\text{rank}(T) = r$. Then the positive eigenvalues of $T'T$ are equal to the positive eigenvalues of TT' .

Proof. Let $\lambda > 0$ be an eigenvalue of $T'T$ with multiplicity h . Since the $n \times n$ matrix $T'T$ is symmetric, we can find an $n \times h$ matrix X , whose columns are orthonormal, satisfying

$$T'TX = \lambda X.$$

Let $Y = TX$ and observe that

$$TT'Y = TT'TX = T(\lambda X) = \lambda TX = \lambda Y,$$

so that λ is also an eigenvalue of TT' . Its multiplicity is also h because

$$\begin{aligned} \text{rank}(Y) &= \text{rank}(TX) = \text{rank}((TX)'TX) \\ &= \text{rank}(X'T'TX) = \text{rank}(\lambda X'X) \\ &= \text{rank}(\lambda I_h) = h, \end{aligned}$$

and so the proof is complete. \square

Example 3.15 In multivariate multiple regression, we have multiple explanatory variables, x_1, \dots, x_k , as in the standard multiple regression model described in Example 2.11, but the response $\mathbf{y} = (y_1, \dots, y_m)'$ is a random vector instead of a random variable. The model is

$$\mathbf{y} = B'\mathbf{x} + \epsilon,$$

where B is $k \times m$, $\mathbf{x} = (x_1, \dots, x_k)'$, and the $m \times 1$ vector ϵ denotes a random error. If we have N observations of the response vector, $\mathbf{y}_1, \dots, \mathbf{y}_N$, and N corresponding explanatory vectors, $\mathbf{x}_1, \dots, \mathbf{x}_N$, the model can be written as

$$Y = XB + E,$$

where $Y' = (\mathbf{y}_1, \dots, \mathbf{y}_N)$, $X' = (\mathbf{x}_1, \dots, \mathbf{x}_N)$, and E is an $N \times m$ matrix containing error terms. The least squares estimator of B is the $k \times m$ matrix which minimizes the sum of squares of the matrix $(Y - XB_0)$, that is, minimizes $\text{tr}\{(Y - XB_0)'(Y - XB_0)\}$, over all choices for B_0 . We now show that this least squares estimator is $\hat{B} = (X'X)^{-1}X'Y$ when X has full column rank. Note that

$$\begin{aligned} (Y - XB_0)'(Y - XB_0) &= (Y - X\hat{B} + X\hat{B} - XB_0)'(Y - X\hat{B} + X\hat{B} - XB_0) \\ &= \{Y - X\hat{B} + X(\hat{B} - B_0)\}'\{Y - X\hat{B} + X(\hat{B} - B_0)\} \\ &= (Y - X\hat{B})'(Y - X\hat{B}) + (\hat{B} - B_0)'X'X(\hat{B} - B_0) \\ &\quad + (Y - X\hat{B})'X(\hat{B} - B_0) + (\hat{B} - B_0)'X'(Y - X\hat{B}) \\ &= (Y - X\hat{B})'(Y - X\hat{B}) + (\hat{B} - B_0)'X'X(\hat{B} - B_0), \end{aligned}$$

since $X'(Y - X\hat{B}) = X'\{Y - X(X'X)^{-1}X'Y\} = (X' - X')Y = (0)$. Thus,

$$\text{tr}\{(Y - XB_0)'(Y - XB_0)\} \geq \text{tr}\{(Y - X\hat{B})'(Y - X\hat{B})\},$$

since $(\hat{B} - B_0)'X'X(\hat{B} - B_0)$ is nonnegative definite.

Next we will use the Courant–Fischer min–max theorem to prove the following important monotonicity property of the eigenvalues of symmetric matrices.

Theorem 3.28 Let A be an $m \times m$ symmetric matrix and B be an $m \times m$ nonnegative definite matrix. Then, for $h = 1, \dots, m$, we have

$$\lambda_h(A + B) \geq \lambda_h(A),$$

where the inequality is strict if B is positive definite.

Proof. For an arbitrary $m \times (h - 1)$ matrix B_h satisfying $B_h' B_h = I_{h-1}$, we have

$$\begin{aligned} \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'(A + B)x}{x'x} &= \max_{\substack{B_h' x=0 \\ x \neq 0}} \left(\frac{x'Ax}{x'x} + \frac{x'Bx}{x'x} \right) \\ &\geq \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} + \min_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Bx}{x'x} \\ &\geq \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} + \min_{x \neq 0} \frac{x'Bx}{x'x} \\ &= \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} + \lambda_m(B) \\ &\geq \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Ax}{x'x}, \end{aligned}$$

where the last equality follows from Theorem 3.16. The final inequality is strict if B is positive definite because, in this case, $\lambda_m(B) > 0$. Now minimizing both sides of the equation above over all choices of B_h satisfying $B_h' B_h = I_{h-1}$ and using (3.9) of Theorem 3.18, we get

$$\begin{aligned} \lambda_h(A + B) &= \min_{B_h} \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'(A + B)x}{x'x} \\ &\geq \min_{B_h} \max_{\substack{B_h' x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} = \lambda_h(A). \end{aligned}$$

This completes the proof. \square

Note that there is not a general bounding relationship between $\lambda_h(A + B)$ and $\lambda_h(A) + \lambda_h(B)$. For instance, if $A = \text{diag}(1, 2, 3, 4)$ and $B = \text{diag}(8, 6, 4, 2)$, then

$$\lambda_2(A + B) = 8 < \lambda_2(A) + \lambda_2(B) = 3 + 6 = 9,$$

whereas

$$\lambda_3(A + B) = 7 > \lambda_3(A) + \lambda_3(B) = 2 + 4 = 6.$$

In Example 3.12, we discussed a situation in which the eigenvalues and eigenvectors of

$$A = \sum_{i=1}^k (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})'$$

were used in analyzing differences among the group means $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$. For instance, an eigenvector \mathbf{x}_1 , corresponding to the largest eigenvalue of A , gives the direction of maximum dispersion among the group means in that

$$\frac{\mathbf{x}_1' A \mathbf{x}_1}{\mathbf{x}_1' \mathbf{x}_1}$$

is maximized. The division here by $\mathbf{x}_1' \mathbf{x}_1$, which removes the effect of scale, may not be appropriate if the groups have covariance matrices other than the identity matrix. Suppose, for example, that each group has the same covariance matrix B . If \mathbf{y} is a random vector with covariance matrix B , then the variability of \mathbf{y} in the direction given by \mathbf{x} will be $\text{var}(\mathbf{x}'\mathbf{y}) = \mathbf{x}' B \mathbf{x}$. Since differences among the groups in a direction with high variability will not be as important as similar differences in another direction with low variability, we will adjust for these differences in variability by constructing the ratio

$$\frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' B \mathbf{x}}.$$

The vector \mathbf{x}_1 that maximizes this ratio will then identify the one-dimensional subspace of R^m in which the group means differ the most, when adjusting for differences in variability. The next step after finding \mathbf{x}_1 would be to find the vector \mathbf{x}_2 that maximizes this ratio but has $\mathbf{x}_2' \mathbf{y}$ uncorrelated with $\mathbf{x}_1' \mathbf{y}$; this would be the vector \mathbf{x}_2 that maximizes the ratio above subject to the constraint that $\mathbf{x}_1' B \mathbf{x}_2 = 0$. Continuing in this fashion, we would determine the m vectors $\mathbf{x}_1, \dots, \mathbf{x}_m$ that yield the m extremal values $\lambda_1, \dots, \lambda_m$ of the ratio. These extremal values are identified in the following theorem.

Theorem 3.29 Let A and B be $m \times m$ matrices, with A being symmetric and B being positive definite. Then the eigenvalues of $B^{-1}A$, $\lambda_1(B^{-1}A) \geq \dots \geq \lambda_m(B^{-1}A)$, are real and there exists a linearly independent set of eigenvectors, $\mathbf{x}_1, \dots, \mathbf{x}_m$, corresponding to these eigenvalues. In addition, if we define $X_h = (\mathbf{x}_1, \dots, \mathbf{x}_h)$ and $Y_h = (\mathbf{x}_h, \dots, \mathbf{x}_m)$ for $h = 1, \dots, m$, then

$$\lambda_h(B^{-1}A) = \min_{\substack{Y_{h+1}' B \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' B \mathbf{x}}$$

and

$$\lambda_h(B^{-1}A) = \max_{\substack{X'_{h-1}B\mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'B\mathbf{x}},$$

where the min and max are over all $\mathbf{x} \neq 0$ when $h = m$ and $h = 1$, respectively.

Proof. Let $B = PDP'$ be the spectral decomposition of B , so that $D = \text{diag}(d_1, \dots, d_m)$, where the eigenvalues of B , d_1, \dots, d_m , are all positive because of Theorem 3.25. If we let $T = PD^{1/2}P'$, where $D^{1/2} = \text{diag}(d_1^{1/2}, \dots, d_m^{1/2})$, then $B = TT = T^2$ and T , like B , is symmetric and nonsingular. Now it follows from Theorem 3.2(d) that the eigenvalues of $B^{-1}A$ are the same as those of $T^{-1}AT^{-1}$, and these must be real because $T^{-1}AT^{-1}$ is symmetric. Also because of its symmetry, $T^{-1}AT^{-1}$ must have an orthonormal set of eigenvectors, $\mathbf{y}_1, \dots, \mathbf{y}_m$. Note that if we write $\lambda_i = \lambda_i(B^{-1}A) = \lambda_i(T^{-1}AT^{-1})$, then $T^{-1}AT^{-1}\mathbf{y}_i = \lambda_i\mathbf{y}_i$, so that

$$T^{-1}T^{-1}AT^{-1}\mathbf{y}_i = \lambda_i T^{-1}\mathbf{y}_i$$

or

$$B^{-1}A(T^{-1}\mathbf{y}_i) = \lambda_i(T^{-1}\mathbf{y}_i).$$

Thus, $\mathbf{x}_i = T^{-1}\mathbf{y}_i$ is an eigenvector of $B^{-1}A$ corresponding to the eigenvalue $\lambda_i = \lambda_i(B^{-1}A)$ and $\mathbf{y}_i = T\mathbf{x}_i$. Clearly the vectors $\mathbf{x}_1, \dots, \mathbf{x}_m$ are linearly independent because $\mathbf{y}_1, \dots, \mathbf{y}_m$ are orthonormal. All that remains is to prove the identities involving the minimum and maximum. We will just prove the result involving the minimum; the proof for the maximum is similar. Putting $\mathbf{y} = T\mathbf{x}$, we find that

$$\begin{aligned} \min_{\substack{Y'_{h+1}B\mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'B\mathbf{x}} &= \min_{\substack{Y'_{h+1}T\mathbf{x}=0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}'TT^{-1}AT^{-1}T\mathbf{x}}{\mathbf{x}'TT\mathbf{x}} \\ &= \min_{\substack{Y'_{h+1}T\mathbf{y}=0 \\ \mathbf{y} \neq 0}} \frac{\mathbf{y}'T^{-1}AT^{-1}\mathbf{y}}{\mathbf{y}'\mathbf{y}}. \end{aligned} \quad (3.24)$$

Since the rows of $Y'_{h+1}T$ are the transposes of the eigenvectors $T\mathbf{x}_{h+1}, \dots, T\mathbf{x}_m$ of $T^{-1}AT^{-1}$, it follows from Theorem 3.17 that (3.24) equals $\lambda_h(T^{-1}AT^{-1})$, which we have already established as being the same as $\lambda_h(B^{-1}A)$. \square

Note that if \mathbf{x}_i is an eigenvector of $B^{-1}A$ corresponding to the eigenvalue $\lambda_i = \lambda_i(B^{-1}A)$, then

$$B^{-1}A\mathbf{x}_i = \lambda_i\mathbf{x}_i$$

or, equivalently,

$$A\mathbf{x}_i = \lambda_i B\mathbf{x}_i. \quad (3.25)$$

Equation (3.25) is similar to the eigenvalue-eigenvector equation of A , except for the multiplication of \mathbf{x}_i by B on the right-hand side of the equation. The eigenvalues satisfying (3.25) are sometimes referred to as the eigenvalues of A in the metric of B . Note that if we premultiply (3.25) by \mathbf{x}'_i and then solve for λ_i , we get

$$\lambda_i(B^{-1}A) = \frac{\mathbf{x}'_i A \mathbf{x}_i}{\mathbf{x}'_i B \mathbf{x}_i};$$

that is, the extremal values given in Theorem 3.29 are attained at the eigenvectors of $B^{-1}A$.

The result given in Theorem 3.29 can be generalized just as the result in Theorem 3.17 was generalized to that given in Theorem 3.18.

Theorem 3.30 Let A and B be $m \times m$ matrices, with A being symmetric and B being positive definite. For $h = 1, \dots, m$, let B_h be any $m \times (h-1)$ matrix and C_h any $m \times (m-h)$ matrix satisfying, $B'_h B_h = I_{h-1}$ and $C'_h C_h = I_{m-h}$. Then

$$\lambda_h(B^{-1}A) = \min_{B_h} \max_{\substack{B'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' B \mathbf{x}},$$

and

$$\lambda_h(B^{-1}A) = \max_{C_h} \min_{\substack{C'_h \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' B \mathbf{x}},$$

where the inner max and min are over all $\mathbf{x} \neq \mathbf{0}$ when $h = 1$ and $h = m$, respectively.

The proof of Theorem 3.29 suggests a way of simultaneously diagonalizing the matrices A and B . Since $T^{-1}AT^{-1}$ is a symmetric matrix, it can be expressed in the form $Q\Lambda Q'$, where Q is an orthogonal matrix and Λ is the diagonal matrix $\text{diag}(\lambda_1(T^{-1}AT^{-1}), \dots, \lambda_m(T^{-1}AT^{-1}))$. The matrix $C = Q'T^{-1}$ is nonsingular because Q and T^{-1} are nonsingular, and

$$\begin{aligned} CAC' &= Q'T^{-1}AT^{-1}Q = Q'Q\Lambda Q'Q = \Lambda, \\ CBC' &= Q'T^{-1}TTT^{-1}Q = Q'Q = I_m. \end{aligned}$$

Equivalently, if $G = C^{-1}$, we have $A = G\Lambda G'$ and $B = GG'$. This simultaneous diagonalization is useful in proving our next result in Theorem 3.31. For some other related results, see Olkin and Tomsy (1981).

Theorem 3.31 Let A be an $m \times m$ symmetric matrix and B be an $m \times m$ positive definite matrix. If F is any $m \times h$ matrix with full column rank, then for $i = 1, \dots, h$

$$\lambda_i((F'BF)^{-1}(F'AF)) \leq \lambda_i(B^{-1}A),$$

and further

$$\max_F \lambda_i((F'BF)^{-1}(F'AF)) = \lambda_i(B^{-1}A).$$

Proof. Note that the second equation implies the first, so our proof simply involves the verification of the second equation. Let the nonsingular $m \times m$ matrix G be such that $B = GG'$ and $A = G\Lambda G'$, where $\Lambda = \text{diag}(\lambda_1(B^{-1}A), \dots, \lambda_m(B^{-1}A))$. Then

$$\begin{aligned} \max_F \lambda_i((F'BF)^{-1}(F'AF)) &= \max_F \lambda_i((F'GG'F)^{-1}(F'G\Lambda G'F)) \\ &= \max_E \lambda_i((E'E)^{-1}(E'\Lambda E)), \end{aligned}$$

where this last maximization is also over all $m \times h$ matrices of rank h , because $E = G'F$ must have the same rank as F . Note that because E has rank h , the $h \times h$ matrix $E'E$ is a nonsingular symmetric matrix. As was seen in the proof of Theorem 3.29, such a matrix can be expressed as $E'E = TT'$ for some nonsingular symmetric $h \times h$ matrix T . It then follows that

$$\begin{aligned} \max_E \lambda_i((E'E)^{-1}(E'\Lambda E)) &= \max_E \lambda_i((TT')^{-1}(E'\Lambda E)) \\ &= \max_E \lambda_i(T^{-1}E'\Lambda ET^{-1}), \end{aligned}$$

where this last equality follows from Theorem 3.2(d). Now if we define the $m \times h$ rank h matrix $H = ET^{-1}$, then

$$H'H = T^{-1}E'E T^{-1} = T^{-1}TTT^{-1} = I_h.$$

Thus,

$$\max_E \lambda_i(T^{-1}E'\Lambda ET^{-1}) = \max_H \lambda_i(H'\Lambda H) = \lambda_i(B^{-1}A),$$

where the final equality follows from Theorem 3.19 and the fact that equality is actually achieved with the choice of $H' = [I_h \quad (0)]$. \square

Example 3.16 Many multivariate analyses are simply generalizations or extensions of corresponding univariate analyses. In this example, we begin with what is known as the univariate one-way classification model in which we have independent samples of a response y from k different populations or treatments, with a sample size of n_i from the i th population. The j th observation from the i th sample can be expressed as

$$y_{ij} = \mu_i + \epsilon_{ij},$$

where the μ_i 's are constants and the ϵ_{ij} 's are independent and identically distributed as $N(0, \sigma^2)$. Our goal is to determine if the μ_i 's are all the same; that is, we wish to test the null hypothesis $H_0 : \mu_1 = \dots = \mu_k$ against the alternative hypothesis H_1 : at

least two of the μ_i 's differ. An analysis of variance compares (see Problem 2.45) the variability between treatments,

$$\text{SST} = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2,$$

to the variability within treatments,

$$\text{SSE} = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2,$$

where

$$\bar{y}_i = \sum_{j=1}^{n_i} y_{ij} / n_i, \quad \bar{y} = \sum_{i=1}^k n_i \bar{y}_i / n, \quad n = \sum_{i=1}^k n_i.$$

SST is referred to as the sum of squares for treatment whereas SSE is called the sum of squares for error. The hypothesis H_0 is rejected if the statistic

$$F = \frac{\text{SST}/(k-1)}{\text{SSE}/(n-k)}$$

exceeds the appropriate quantile of the F distribution with $k-1$ and $n-k$ degrees of freedom. Now suppose that instead of obtaining the value of one response variable for each observation, we obtain the values of m different response variables for each observation. If \mathbf{y}_{ij} is the $m \times 1$ vector of responses obtained as the j th observation from the i th treatment, then we have the multivariate one-way classification model given by

$$\mathbf{y}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_{ij},$$

where $\boldsymbol{\mu}_i$ is an $m \times 1$ vector of constants and $\boldsymbol{\epsilon}_{ij} \sim N_m(\mathbf{0}, \Omega)$, independently. Measures of the between treatment variability and within treatment variability are now given by the matrices,

$$B = \sum_{i=1}^k n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)'.$$

One approach to testing the null hypothesis $H_0 : \boldsymbol{\mu}_1 = \cdots = \boldsymbol{\mu}_k$ against the alternative hypothesis H_1 : at least two of the $\boldsymbol{\mu}_i$'s differ, is by a method called the union-intersection procedure. This technique is based on the following decomposition of the hypotheses H_0 and H_1 into univariate hypotheses. If \mathbf{c} is any $m \times 1$ vector, and we define the hypothesis $H_0(\mathbf{c}) : \mathbf{c}'\boldsymbol{\mu}_1 = \cdots = \mathbf{c}'\boldsymbol{\mu}_k$, then the intersection of $H_0(\mathbf{c})$ over all $\mathbf{c} \in R^m$ is the hypothesis H_0 . In addition, if we define the hypothesis $H_1(\mathbf{c})$: at least two of the $\mathbf{c}'\boldsymbol{\mu}_i$'s differ, then the union of the hypotheses $H_1(\mathbf{c})$ over all $\mathbf{c} \in R^m$ is the hypothesis H_1 . Thus, we should reject

the hypothesis H_0 if and only if we reject $H_0(\mathbf{c})$ for at least one \mathbf{c} . Now the null hypothesis $H_0(\mathbf{c})$ involves the univariate one-way classification model in which $\mathbf{c}'\mathbf{y}_{ij}$ is the response, and so we would reject $H_0(\mathbf{c})$ for large values of the F statistic

$$F(\mathbf{c}) = \frac{\text{SST}(\mathbf{c})/(k-1)}{\text{SSE}(\mathbf{c})/(n-k)},$$

where $\text{SST}(\mathbf{c})$ and $\text{SSE}(\mathbf{c})$ are the sums of squares for treatments and errors, respectively, computed for the responses $\mathbf{c}'\mathbf{y}_{ij}$. Since H_0 is rejected if $H_0(\mathbf{c})$ is rejected for at least one \mathbf{c} , we will reject H_0 if $F(\mathbf{c})$ is sufficiently large for at least one \mathbf{c} or, equivalently, if

$$\max_{\mathbf{c} \neq \mathbf{0}} F(\mathbf{c})$$

is sufficiently large. Omitting the constants $(k-1)$ and $(n-k)$ and noting that the sums of squares $\text{SST}(\mathbf{c})$ and $\text{SSE}(\mathbf{c})$ can be expressed using B and W as

$$\text{SST}(\mathbf{c}) = \mathbf{c}'B\mathbf{c}, \quad \text{SSE}(\mathbf{c}) = \mathbf{c}'W\mathbf{c},$$

we find that we reject H_0 for large values of

$$\max_{\mathbf{c} \neq \mathbf{0}} \frac{\mathbf{c}'B\mathbf{c}}{\mathbf{c}'W\mathbf{c}} = \lambda_1(W^{-1}B), \quad (3.26)$$

where the right-hand side follows from Theorem 3.29. Thus, if $u_{1-\alpha}$ is the $(1-\alpha)$ th quantile of the distribution of the largest eigenvalue $\lambda_1(W^{-1}B)$ (see, for example, Morrison, 2005) so that

$$P[\lambda_1(W^{-1}B) \leq u_{1-\alpha} | H_0] = 1 - \alpha, \quad (3.27)$$

then we would reject H_0 if $\lambda_1(W^{-1}B) > u_{1-\alpha}$. One advantage of the union-intersection procedure is that it naturally leads to simultaneous confidence intervals. It follows immediately from (3.26) and (3.27) that for any mean vectors $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$, with probability $1 - \alpha$, the inequality

$$\frac{\sum_{i=1}^k n_i \mathbf{c}' \{(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}) - (\boldsymbol{\mu}_i - \boldsymbol{\mu})\} \{(\bar{\mathbf{y}}_i - \bar{\mathbf{y}}) - (\boldsymbol{\mu}_i - \boldsymbol{\mu})\}' \mathbf{c}}{\mathbf{c}'W\mathbf{c}} \leq u_{1-\alpha}, \quad (3.28)$$

holds for all $m \times 1$ vectors \mathbf{c} , where

$$\boldsymbol{\mu} = \sum_{i=1}^k n_i \boldsymbol{\mu}_i / n.$$

Scheffé's method (see, Scheffé, 1953, or Miller, 1981) can then be used on (3.28) to yield the inequalities

$$\begin{aligned}
 \sum_{i=1}^k \sum_{j=1}^m a_i c_j \bar{y}_{ij} - \sqrt{u_{1-\alpha} \mathbf{c}' W \mathbf{c} \left(\sum_{i=1}^k a_i^2 / n_i \right)} \\
 \leq \sum_{i=1}^k \sum_{j=1}^m a_i c_j \mu_{ij} \\
 \leq \sum_{i=1}^k \sum_{j=1}^m a_i c_j \bar{y}_{ij} + \sqrt{u_{1-\alpha} \mathbf{c}' W \mathbf{c} \left(\sum_{i=1}^k a_i^2 / n_i \right)},
 \end{aligned}$$

which hold with probability $1 - \alpha$, for all $m \times 1$ vectors \mathbf{c} and all $k \times 1$ vectors \mathbf{a} satisfying $\mathbf{a}' \mathbf{1}_k = 0$.

The remaining results in this section relate the eigenvalues of a matrix product to products of the eigenvalues of the individual matrices. Theorem 3.32, which is due to Anderson and Das Gupta (1963), gives bounds for a single eigenvalue of a matrix product.

Theorem 3.32 Let A be an $m \times m$ nonnegative definite matrix, and let B be an $m \times m$ positive definite matrix. If i, j , and k are integers between 1 and m inclusive and satisfying $j + k \leq i + 1$, then

- (a) $\lambda_i(AB) \leq \lambda_j(A) \lambda_k(B)$,
- (b) $\lambda_{m-i+1}(AB) \geq \lambda_{m-j+1}(A) \lambda_{m-k+1}(B)$.

Proof. Let the columns of the $m \times (j - 1)$ matrix H_1 be orthonormal eigenvectors of A corresponding to $\lambda_1(A), \dots, \lambda_{j-1}(A)$, and let the columns of the $m \times (k - 1)$ matrix H_2 be orthonormal eigenvectors of B corresponding to $\lambda_1(B), \dots, \lambda_{k-1}(B)$. Define the $m \times (j + k - 2)$ matrix H as $H = [H_1 \ H_2]$. Then

$$\begin{aligned}
 \lambda_i(AB) &\leq \lambda_{j+k-1}(AB) \\
 &\leq \max_{\substack{H'x=0 \\ x \neq 0}} \frac{x'Ax}{x'B^{-1}x} \\
 &= \max_{\substack{H'x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} \frac{x'x}{x'B^{-1}x} \\
 &\leq \max_{\substack{H'x=0 \\ x \neq 0}} \frac{x'Ax}{x'x} \max_{\substack{H'x=0 \\ x \neq 0}} \frac{x'x}{x'B^{-1}x}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \max_{\substack{H'_1 \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \max_{\substack{H'_2 \mathbf{x} = \mathbf{0} \\ \mathbf{x} \neq \mathbf{0}}} \frac{\mathbf{x}' \mathbf{x}}{\mathbf{x}' B^{-1} \mathbf{x}} \\
 &= \lambda_j(A) \lambda_k(B),
 \end{aligned}$$

where the second inequality follows from Theorem 3.30 and the final equality follows from Theorem 3.17. This establishes the inequality given in (a). The inequality in (b) can be obtained from (a) as follows. Write $A_* = TAT$, where T is the symmetric matrix satisfying $B = T^2$. Then A_* is nonnegative definite because A is. Applying (a) to A_* and B^{-1} and using the fact that $\lambda_i(A_* B^{-1}) = \lambda_i(A)$ and $\lambda_j(A_*) = \lambda_j(AB)$, we get

$$\lambda_i(A) \leq \lambda_j(AB) \lambda_k(B^{-1}).$$

Since $\lambda_k(B^{-1}) = \lambda_{m-k+1}^{-1}(B)$, this leads to

$$\lambda_j(AB) \geq \lambda_i(A) \lambda_{m-k+1}(B). \quad (3.29)$$

For each (i, j, k) satisfying the constraint given in the theorem so also will (i_*, j_*, k) , where $i_* = m - j + 1$ and $j_* = m - i + 1$. Making these substitutions in (3.29) yields (b). \square

Our next result, due to Lidskiĭ (1950), gives a bound for a product of eigenvalues of a matrix product. For a proof of this result, see Zhang (2011).

Theorem 3.33 Let A and B be $m \times m$ nonnegative definite matrices. If i_1, \dots, i_k are integers satisfying $1 \leq i_1 < \dots < i_k \leq m$, then

$$\prod_{j=1}^k \lambda_{i_j}(AB) \leq \prod_{j=1}^k \lambda_{i_j}(A) \lambda_j(B),$$

for $k = 1, \dots, m$, with equality for $k = m$.

Some additional inequalities for eigenvalues will be given in Sections 7.6 and 10.2.

3.9 ANTIEIGENVALUES AND ANTIEIGENVECTORS

Consider an $m \times m$ positive definite matrix A and an $m \times 1$ nonnull vector \mathbf{x} . If $\mathbf{y} = A\mathbf{x}$, it follows from the Cauchy-Schwarz inequality that

$$\mathbf{x}' A \mathbf{x} = \mathbf{x}' \mathbf{y} \leq \sqrt{(\mathbf{x}' \mathbf{x})(\mathbf{y}' \mathbf{y})} = \sqrt{(\mathbf{x}' \mathbf{x})(\mathbf{x}' A^2 \mathbf{x})},$$

with equality if and only if one of the vectors, \mathbf{x} and \mathbf{y} , is a scalar multiple of the other, that is, $A\mathbf{x} = \lambda\mathbf{x}$ for some scalar λ . Another way of stating this is that the function

$$\psi(\mathbf{x}) = \frac{\mathbf{x}'A\mathbf{x}}{\sqrt{(\mathbf{x}'\mathbf{x})(\mathbf{x}'A^2\mathbf{x})}}$$

has a maximum value of 1, which is attained if and only if \mathbf{x} is an eigenvector of A . This is not surprising since $\psi(\mathbf{x})$ is $\cos \theta$, where θ is the angle between \mathbf{x} and $A\mathbf{x}$, and clearly the eigenvectors of A minimize this angle.

While the eigenvectors of A maximize $\psi(\mathbf{x})$, the vectors that minimize $\psi(\mathbf{x})$ are known as the antieigenvectors of A . The notion of antieigenvalues and antieigenvectors was originated by Gustafson (1972).

Definition 3.2 Let A be an $m \times m$ positive definite matrix and \mathbf{x} be an $m \times 1$ vector. Then

$$\mu_1 = \min_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}'A\mathbf{x}}{\sqrt{(\mathbf{x}'\mathbf{x})(\mathbf{x}'A^2\mathbf{x})}}$$

is the first antieigenvalue of A and \mathbf{x} is a corresponding antieigenvector if $\mu_1 = \psi(\mathbf{x})$.

The quantity $\theta = \cos^{-1}(\mu_1)$ can be described as the largest turning angle of A since it gives the maximum angle between \mathbf{x} and $A\mathbf{x}$ over all choices of $\mathbf{x} \neq \mathbf{0}$.

Let $\lambda_1 \geq \dots \geq \lambda_m > 0$ be the eigenvalues of A and $\mathbf{x}_1, \dots, \mathbf{x}_m$ corresponding orthonormal eigenvectors. The following result shows that the first antieigenvalue and corresponding normalized first antieigenvectors can be expressed in terms of $\lambda_1, \lambda_m, \mathbf{x}_1$, and \mathbf{x}_m .

Theorem 3.34 The $m \times m$ positive definite matrix A has first antieigenvalue

$$\mu_1 = \frac{2\sqrt{\lambda_1\lambda_m}}{\lambda_1 + \lambda_m},$$

with corresponding normalized antieigenvectors given by

$$\begin{aligned} \mathbf{x}_{1+} &= \left(\frac{\lambda_m}{\lambda_1 + \lambda_m} \right)^{1/2} \mathbf{x}_1 + \left(\frac{\lambda_1}{\lambda_1 + \lambda_m} \right)^{1/2} \mathbf{x}_m, \\ \mathbf{x}_{1-} &= \left(\frac{\lambda_m}{\lambda_1 + \lambda_m} \right)^{1/2} \mathbf{x}_1 - \left(\frac{\lambda_1}{\lambda_1 + \lambda_m} \right)^{1/2} \mathbf{x}_m. \end{aligned}$$

Proof. Let $A = X\Lambda X'$ be the spectral decomposition of A , and define $\mathbf{y} = \Lambda^{1/2}X'\mathbf{x}$, $\mathbf{z} = (\mathbf{y}'\mathbf{y})^{-1/2}\mathbf{y}$, and $\gamma_i = \lambda_i/(\lambda_1\lambda_m)^{1/2}$. Then

$$\begin{aligned}
 \frac{\mathbf{x}'A\mathbf{x}}{\sqrt{(\mathbf{x}'\mathbf{x})(\mathbf{x}'A^2\mathbf{x})}} &= \frac{(\mathbf{x}'X\Lambda^{1/2})(\Lambda^{1/2}X'\mathbf{x})}{\sqrt{(\mathbf{x}'X\Lambda^{1/2})\Lambda^{-1}(\Lambda^{1/2}X'\mathbf{x})(\mathbf{x}'X\Lambda^{1/2})\Lambda(\Lambda^{1/2}X'\mathbf{x})}} \\
 &= \frac{\mathbf{y}'\mathbf{y}}{\sqrt{(\mathbf{y}'\Lambda^{-1}\mathbf{y})(\mathbf{y}'\Lambda\mathbf{y})}} \\
 &= \frac{1}{\sqrt{(\mathbf{z}'\Lambda^{-1}\mathbf{z})(\mathbf{z}'\Lambda\mathbf{z})}} \\
 &= \frac{1}{\sqrt{\left(\sum_{i=1}^m \lambda_i^{-1}z_i^2\right)\left(\sum_{i=1}^m \lambda_i z_i^2\right)}} \\
 &= \frac{1}{\sqrt{\left(\sum_{i=1}^m \gamma_i^{-1}z_i^2\right)\left(\sum_{i=1}^m \gamma_i z_i^2\right)}}. \tag{3.30}
 \end{aligned}$$

Now since the function $g(\gamma) = \gamma + \gamma^{-1}$ has positive second derivative and $g(\gamma_1) = g(\gamma_m)$, it follows that $g(\gamma_j) \leq g(\gamma_1)$ for $j = 2, \dots, m-1$, and so

$$\begin{aligned}
 \frac{1}{2} \left(\sum_{i=1}^m \gamma_i z_i^2 + \sum_{i=1}^m \gamma_i^{-1} z_i^2 \right) &= \frac{1}{2} \sum_{i=1}^m (\gamma_i + \gamma_i^{-1}) z_i^2 \\
 &\leq \frac{1}{2} \sum_{i=1}^m (\gamma_1 + \gamma_1^{-1}) z_i^2 \\
 &= \frac{1}{2} (\gamma_1 + \gamma_1^{-1}) = \frac{1}{2} (\gamma_1 + \gamma_m). \tag{3.31}
 \end{aligned}$$

An application of the arithmetic-geometric mean inequality (see Section 10.6) yields

$$\frac{1}{2} \left(\sum_{i=1}^m \gamma_i z_i^2 + \sum_{i=1}^m \gamma_i^{-1} z_i^2 \right) \geq \left\{ \left(\sum_{i=1}^m \gamma_i z_i^2 \right) \left(\sum_{i=1}^m \gamma_i^{-1} z_i^2 \right) \right\}^{1/2}. \tag{3.32}$$

Combining (3.31) and (3.32), we have

$$\left\{ \left(\sum_{i=1}^m \gamma_i z_i^2 \right) \left(\sum_{i=1}^m \gamma_i^{-1} z_i^2 \right) \right\}^{1/2} \leq \frac{1}{2} (\gamma_1 + \gamma_m) = \frac{1}{2} \left(\frac{\lambda_1 + \lambda_m}{\sqrt{\lambda_1 \lambda_m}} \right) = \mu_1^{-1},$$

which establishes μ_1 as a lower bound for (3.30). This bound is attained if we have equality in (3.31) and (3.32), so that $z_1^2 = z_m^2 = \frac{1}{2}$ or $y_1 = y_m$, $y_2 = \dots = y_{m-1} = 0$.

Consequently, the bound is attained when

$$\mathbf{x} = X\Lambda^{-1/2}\mathbf{y} = \frac{y_1}{\sqrt{\lambda_1}}\mathbf{x}_1 + \frac{y_1}{\sqrt{\lambda_m}}\mathbf{x}_m.$$

If $\lambda_1 > \lambda_2$ and $\lambda_{m-1} > \lambda_m$, there are only two linearly independent vectors of this form, which, when normalized, can be expressed as \mathbf{x}_{1+} and \mathbf{x}_{1-} . \square

Additional antieigenvalues and associated antieigenvectors can be defined as follows. Let $X_k = (\mathbf{x}_1, \dots, \mathbf{x}_{k-1}, \mathbf{x}_{m-k+2}, \dots, \mathbf{x}_m)$. Then for $k = 2, \dots, r$, where $r = m/2$ if m is even and $r = (m-1)/2$ if m is odd, the k th antieigenvalue of A is

$$\mu_k = \min_{\substack{\mathbf{x}'_k \mathbf{x} = 0 \\ \mathbf{x} \neq 0}} \frac{\mathbf{x}' A \mathbf{x}}{\sqrt{(\mathbf{x}' \mathbf{x})(\mathbf{x}' A^2 \mathbf{x})}},$$

with associated antieigenvectors given by any \mathbf{x} satisfying $\psi(\mathbf{x}) = \mu_k$. The proof of the following result, which is omitted, is similar to the proof of Theorem 3.34.

Theorem 3.35 The $m \times m$ positive definite matrix A has k th antieigenvalue

$$\mu_k = \frac{2\sqrt{\lambda_k \lambda_{m-k+1}}}{\lambda_k + \lambda_{m-k+1}},$$

for $k = 2, \dots, r$. Corresponding normalized antieigenvectors are given by

$$\begin{aligned} \mathbf{x}_{k+} &= \left(\frac{\lambda_{m-k+1}}{\lambda_k + \lambda_{m-k+1}} \right)^{1/2} \mathbf{x}_k + \left(\frac{\lambda_k}{\lambda_k + \lambda_{m-k+1}} \right)^{1/2} \mathbf{x}_{m-k+1}, \\ \mathbf{x}_{k-} &= \left(\frac{\lambda_{m-k+1}}{\lambda_k + \lambda_{m-k+1}} \right)^{1/2} \mathbf{x}_k - \left(\frac{\lambda_k}{\lambda_k + \lambda_{m-k+1}} \right)^{1/2} \mathbf{x}_{m-k+1}. \end{aligned}$$

Applications of antieigenvalues and antieigenvectors in statistics can be found in Khattree (2003), Rao (2005), and Gustafson (2006).

PROBLEMS

3.1 Consider the 3×3 matrix

$$A = \begin{bmatrix} 9 & -3 & -4 \\ 12 & -4 & -6 \\ 8 & -3 & -3 \end{bmatrix}.$$

(a) Find the eigenvalues of A .

- (b) Find a normalized eigenvector corresponding to each eigenvalue.
 (c) Find $\text{tr}(A^{10})$.
- 3.2** Consider the 3×3 matrix A given in Problem 3.1 and the vector $\mathbf{x} = (3, 3, 4)'$.
 (a) Show that \mathbf{x} can be written as a linear combination of the eigenvectors of A .
 (b) Find $A^{10}\mathbf{x}$ without actually computing A^{10} .
- 3.3** Find the eigenvalues of A' , where A is the matrix given in Problem 3.1. Determine the eigenspaces for A' , and compare these with those of A .
- 3.4** Let the 3×3 matrix A be given by

$$A = \begin{bmatrix} 1 & -2 & 0 \\ 1 & 4 & 0 \\ 0 & 0 & 2 \end{bmatrix}.$$

- (a) Find the eigenvalues of A .
 (b) For each different value of λ , determine the associated eigenspace $S_A(\lambda)$.
 (c) Describe the eigenspaces obtained in part (b).
- 3.5** Consider the 4×4 matrix

$$A = \begin{bmatrix} 0 & 0 & 2 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & -2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}.$$

- (a) Find the eigenvalues of A .
 (b) Find the eigenspaces of A .
- 3.6** If the $m \times m$ matrix A has eigenvalues $\lambda_1, \dots, \lambda_m$ and corresponding eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$, show that the matrix $(A + \gamma I_m)$ has eigenvalues $\lambda_1 + \gamma, \dots, \lambda_m + \gamma$ and corresponding eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$.
- 3.7** In Example 3.7, we discussed the use of principal components regression as a way of overcoming the difficulties associated with multicollinearity. Another approach, called ridge regression, replaces the ordinary least squares estimator in the standardized model, $\hat{\boldsymbol{\delta}}_1 = (Z_1'Z_1)^{-1}Z_1'\mathbf{y}$ by $\hat{\boldsymbol{\delta}}_{1\gamma} = (Z_1'Z_1 + \gamma I_k)^{-1}Z_1'\mathbf{y}$, where γ is a small positive number. This adjustment will reduce the impact of the near singularity of $Z_1'Z_1$ because the addition of γI_k increases each of the eigenvalues of $Z_1'Z_1$ by γ .
- (a) Show that if $N \geq 2k + 1$, there is an $N \times k$ matrix W , such that $\hat{\boldsymbol{\delta}}_{1\gamma}$ is the ordinary least squares estimate of $\boldsymbol{\delta}_1$ in the model

$$\mathbf{y} = \delta_0 \mathbf{1}_N + (Z_1 + W)\boldsymbol{\delta}_1 + \boldsymbol{\epsilon};$$

that is, $\hat{\delta}_{1\gamma}$ can be viewed as the ordinary least squares estimator of δ_1 after we have perturbed the matrix of values for the explanatory variables Z_1 by W .

- (b) Show that a $k \times k$ matrix U exists, such that $\hat{\delta}_{1\gamma}$ is the ordinary least squares estimate of δ_1 in the model

$$\begin{bmatrix} y \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \delta_0 \mathbf{1}_N \\ \mathbf{0} \end{bmatrix} + \begin{bmatrix} Z_1 \\ U \end{bmatrix} \delta_1 + \begin{bmatrix} \epsilon \\ \epsilon_* \end{bmatrix},$$

where $\mathbf{0}$ is a $k \times 1$ vector of zeros and $\epsilon_* \sim N_k(\mathbf{0}, \sigma^2 I_k)$, independently of ϵ . Thus, the ridge regression estimator also can be viewed as the least squares estimator obtained after adding k observations, each having zero for the response variable and the small values in U as the values for the explanatory variables.

3.8 Refer to Example 3.7 and the previous exercise.

- (a) Find the expected values of the principal components regression estimator, $\hat{\delta}_{1*}$, and the ridge regression estimator, $\hat{\delta}_{1\gamma}$, thereby showing that each is a biased estimator of δ_1 .
- (b) Find the covariance matrix of $\hat{\delta}_{1*}$ and show that $\text{var}(\hat{\delta}_1) - \text{var}(\hat{\delta}_{1*})$ is a nonnegative definite matrix, where $\hat{\delta}_1$ is the ordinary least squares estimator of δ_1 .
- (c) Find the covariance matrix of $\hat{\delta}_{1\gamma}$ and show that $\text{tr}\{\text{var}(\hat{\delta}_1) - \text{var}(\hat{\delta}_{1\gamma})\}$ is nonnegative.

3.9 If A and B are $m \times m$ matrices and at least one of them is nonsingular, show that the eigenvalues of AB and BA are the same.

3.10 If λ is a real eigenvalue of the $m \times m$ real matrix A , show that there exist real eigenvectors of A corresponding to the eigenvalue λ .

3.11 For some angle θ , consider the 2×2 matrix

$$P = \begin{bmatrix} \cos \theta & -\sin \theta \\ \sin \theta & \cos \theta \end{bmatrix}.$$

- (a) Show that P is an orthogonal matrix.
- (b) Find the eigenvalues of P .

3.12 Suppose that A is $m \times n$ and B is $n \times m$. Show that if $AB\mathbf{x} = \lambda\mathbf{x}$, where $\lambda \neq 0$ and $\mathbf{x} \neq \mathbf{0}$, then $B\mathbf{x}$ is an eigenvector of BA corresponding to λ . Thus, show that AB and BA have the same nonzero eigenvalues by showing that the number of linearly independent eigenvectors of AB corresponding to $\lambda \neq 0$ is the same as the number of linearly independent eigenvectors of BA corresponding to λ .

3.13 Prove the results given in Theorem 3.2.

3.14 Suppose A is an $m \times m$ skew symmetric matrix. Show that each eigenvalue of A is zero or a pure imaginary number; that is, each eigenvalue is of the form $0 + bi$ for some scalar b .

- 3.15** We know from Theorem 3.2(d) that if $m \times m$ matrices A and C satisfy $C = BAB^{-1}$ for some nonsingular matrix B , then A and C have the same eigenvalues. Show by example that the converse is not true; that is, find matrices A and C that have the same eigenvalues, but do not satisfy $C = BAB^{-1}$ for any B .
- 3.16** Suppose that λ is a simple eigenvalue of the $m \times m$ matrix A . Show that $\text{rank}(A - \lambda I_m) = m - 1$.
- 3.17** If A is an $m \times m$ matrix and $\text{rank}(A - \lambda I_m) = m - 1$, show that λ is an eigenvalue of A with multiplicity of at least one.
- 3.18** Let A be an $m \times m$ matrix.
- (a) Show that if A is nonnegative definite, then A^2 is also nonnegative definite.
 - (b) Show that if A is positive definite, then A^{-1} is positive definite.
- 3.19** Consider the $m \times m$ matrix

$$A = \begin{bmatrix} 1 & 1 & 0 & \cdots & 0 \\ 0 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & \cdots & 1 \end{bmatrix},$$

which has each element on and directly above the diagonal equal to 1. Find the eigenvalues and eigenvectors of A .

- 3.20** Let \mathbf{x} and \mathbf{y} be $m \times 1$ vectors.
- (a) Find the eigenvalues and eigenvectors of the matrix \mathbf{xy}' .
 - (b) Show that if $c = 1 + \mathbf{x}'\mathbf{y} \neq 0$, then $I_m + \mathbf{xy}'$ has an inverse and $(I_m + \mathbf{xy}')^{-1} = I_m - c^{-1}\mathbf{xy}'$.
- 3.21** Let A be an $m \times m$ nonsingular matrix with eigenvalues $\lambda_1, \dots, \lambda_m$ and corresponding eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$. If $I_m + A$ is nonsingular, find the eigenvalues and eigenvectors of
- (a) $(I_m + A)^{-1}$,
 - (b) $A + A^{-1}$,
 - (c) $I_m + A^{-1}$.
- 3.22** Suppose that A is an $m \times m$ nonsingular matrix and the sum of the elements in each row of A is 1. Show that the row sums of A^{-1} are also 1.
- 3.23** Let the $m \times m$ nonsingular matrix A be such that $I_m + A$ is nonsingular, and define

$$B = (I_m + A)^{-1} + (I_m + A^{-1})^{-1}.$$

- (a) Show that if \mathbf{x} is an eigenvector of A corresponding to the eigenvalue λ , then \mathbf{x} is an eigenvector of B corresponding to the eigenvalue 1.
- (b) Use Theorem 1.9 to show that $B = I_m$.

3.24 Consider the 2×2 matrix

$$A = \begin{bmatrix} 4 & 2 \\ 3 & 5 \end{bmatrix}.$$

- (a) Find the characteristic equation of A .
- (b) Illustrate Theorem 3.8 by substituting A for λ in the characteristic equation obtained in (a) and then showing the resulting matrix is the null matrix.
- (c) Rearrange the matrix polynomial equation in (b) to obtain an expression for A^2 as a linear combination of A and I_m .
- (d) In a similar fashion, write A^3 and A^{-1} as linear combinations of A and I_m .

3.25 Consider the general 2×2 matrix

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}.$$

- (a) Find the characteristic equation of A .
 - (b) Obtain expressions for the two eigenvalues of A in terms of the elements of A .
 - (c) When will these eigenvalues be real?
- 3.26** If A is $m \times m$, then a nonnull $m \times 1$ vector \mathbf{x} satisfying $A\mathbf{x} = \lambda\mathbf{x}$ for some scalar λ is more precisely referred to as a right eigenvector of A corresponding to λ . An $m \times 1$ vector \mathbf{y} satisfying $\mathbf{y}'A = \mu\mathbf{y}'$ is referred to as a left eigenvector of A corresponding to μ . Show that if $\lambda \neq \mu$, then \mathbf{x} is orthogonal to \mathbf{y} .
- 3.27** Find the eigenvalues and eigenvectors of the matrix $\mathbf{1}_m \mathbf{1}_m'$.
- 3.28** A 3×3 matrix A has eigenvalues 1, 2, and 3, and corresponding eigenvectors $(1, 1, 1)'$, $(1, 2, 0)'$, and $(2, -1, 6)'$. Find A .
- 3.29** Consider the $m \times m$ matrix $A = \alpha I_m + \beta \mathbf{1}_m \mathbf{1}_m'$, where α and β are scalars.
- (a) Find the eigenvalues and eigenvectors of A .
 - (b) Determine the eigenspaces and associated eigenprojections of A .
 - (c) For which values of α and β will A be nonsingular?
 - (d) Using (a), show that when A is nonsingular, then

$$A^{-1} = \alpha^{-1} I_m - \frac{\beta}{\alpha(\alpha + m\beta)} \mathbf{1}_m \mathbf{1}_m'.$$

- (e) Show that the determinant of A is $\alpha^{m-1}(\alpha + m\beta)$.
- 3.30** Consider the $m \times m$ matrix $A = \alpha I_m + \beta \mathbf{c} \mathbf{c}'$, where α and β are scalars and $\mathbf{c} \neq \mathbf{0}$ is an $m \times 1$ vector.
- (a) Find the eigenvalues and eigenvectors of A .
 - (b) Find the determinant of A .

- (c) Give conditions for A to be nonsingular and find an expression for the inverse of A .

3.31 Let A be the 3×3 matrix given by

$$A = \begin{bmatrix} 2 & -1 & 0 \\ -1 & 1 & 1 \\ 0 & 1 & 2 \end{bmatrix}.$$

- (a) Find the eigenvalues and associated normalized eigenvectors of A .
 (b) What is the rank of A ?
 (c) Find the eigenspaces and associated eigenprojections of A .
 (d) Find $\text{tr}(A^4)$.
- 3.32** Construct a 3×3 symmetric matrix having eigenvalues 18, 21, and 28, and corresponding eigenvectors $(1, 1, 2)'$, $(4, -2, -1)'$, and $(1, 3, -2)'$.
- 3.33** Show that if A is an $m \times m$ symmetric matrix with eigenvalues $\lambda_1, \dots, \lambda_m$, then

$$\sum_{i=1}^m \sum_{j=1}^m a_{ij}^2 = \sum_{i=1}^m \lambda_i^2.$$

- 3.34** Show that the matrix $A = (1 - \rho)I_m + \rho \mathbf{1}_m \mathbf{1}_m'$ is positive definite if and only if $-(m-1)^{-1} < \rho < 1$.
- 3.35** Show that if A is an $m \times m$ symmetric matrix with its eigenvalues equal to its diagonal elements, then A must be a diagonal matrix.
- 3.36** Show that the converse of Theorem 3.28 is not true; that is, find symmetric matrices A and B for which $\lambda_i(A+B) \geq \lambda_i(A)$ for $i = 1, \dots, m$ yet B is not nonnegative definite.
- 3.37** Let A be an $m \times n$ matrix with $\text{rank}(A) = r$. Use the spectral decomposition of $A'A$ to show that an $n \times (n-r)$ matrix X exists, such that

$$AX = (0) \quad \text{and} \quad X'X = I_{n-r}.$$

In a similar fashion, show that an $(m-r) \times m$ matrix Y exists, such that

$$YA = (0) \quad \text{and} \quad YY' = I_{m-r}.$$

3.38 Let A be the 2×3 matrix given by

$$A = \begin{bmatrix} 6 & 4 & 4 \\ 3 & 2 & 2 \end{bmatrix}.$$

Find matrices X and Y satisfying the conditions given in the previous exercise.

- 3.39** An $m \times m$ matrix A is said to be nilpotent if $A^k = (0)$ for some positive integer k .

- (a) Show that all of the eigenvalues of a nilpotent matrix are equal to 0.
 (b) Find a matrix, other than the null matrix, that is nilpotent.

3.40 Complete the details of Example 3.11 by showing that

$$P_{1,n} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad P_{2,n} \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

as $n \rightarrow \infty$.

3.41 Prove Corollary 3.18.2.

3.42 Let A be an $m \times m$ symmetric matrix with eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$, and suppose S is a k -dimensional subspace of R^m .

(a) Show that if $\mathbf{x}'A\mathbf{x}/\mathbf{x}'\mathbf{x} \geq b$ for all $\mathbf{x} \in S$, then $\lambda_k \geq b$.

(b) Show that if $\mathbf{x}'A\mathbf{x}/\mathbf{x}'\mathbf{x} \leq b$ for all $\mathbf{x} \in S$, then $\lambda_{m-k+1} \leq b$.

3.43 Show by example that Theorem 3.21 need not hold if the matrices A and B are not symmetric.

3.44 Prove Theorem 3.20.

3.45 Our proof of Theorem 3.28 utilized (3.9) of Theorem 3.18. Obtain an alternative proof of Theorem 3.28 by using (3.10) of Theorem 3.18.

3.46 Let A be a symmetric matrix with $\lambda_1(A) > 0$. Show that

$$\lambda_1(A) = \max_{\mathbf{x}'A\mathbf{x}=1} \frac{1}{\mathbf{x}'\mathbf{x}}.$$

3.47 Let A be an $m \times m$ symmetric matrix and B be an $m \times m$ positive definite matrix. If F is any $m \times h$ matrix with full column rank, then show the following:

(a) $\lambda_{h-i+1}((F'BF)^{-1}(F'AF)) \geq \lambda_{m-i+1}(B^{-1}A)$, for $i = 1, \dots, h$.

(b) $\min_F \lambda_1((F'BF)^{-1}(F'AF)) = \lambda_{m-h+1}(B^{-1}A)$.

(c) $\min_F \lambda_h((F'BF)^{-1}(F'AF)) = \lambda_m(B^{-1}A)$.

3.48 Suppose A is an $m \times m$ matrix with eigenvalues $\lambda_1, \dots, \lambda_m$ and associated eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$, whereas B is $n \times n$ with eigenvalues $\gamma_1, \dots, \gamma_n$ and eigenvectors $\mathbf{y}_1, \dots, \mathbf{y}_n$. What are the eigenvalues and eigenvectors of the $(m+n) \times (m+n)$ matrix

$$C = \begin{bmatrix} A & (0) \\ (0) & B \end{bmatrix}?$$

Generalize this result by giving the eigenvalues and eigenvectors of the matrix

$$C = \begin{bmatrix} C_1 & (0) & \cdots & (0) \\ (0) & C_2 & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & C_r \end{bmatrix}$$

in terms of the eigenvalues and eigenvectors of the square matrices C_1, \dots, C_r .

3.49 Let

$$T = \begin{bmatrix} 1 & -1 & 2 \\ 2 & 1 & 1 \end{bmatrix}.$$

- (a) Find the eigenvalues and corresponding eigenvectors of TT' .
 - (b) Find the eigenvalues and corresponding eigenvectors of $T'T$.
- 3.50** Let A be an $m \times m$ symmetric matrix. Show that if k is a positive integer, then A^{2k} is nonnegative definite.
- 3.51** Show that if A is a nonnegative definite matrix and $a_{ii} = 0$ for some i , then $a_{ij} = a_{ji} = 0$ for all j .
- 3.52** Let A be an $m \times m$ positive definite matrix and B be an $m \times m$ nonnegative definite matrix.
- (a) Use the spectral decomposition of A to show that

$$|A + B| \geq |A|,$$

with equality if and only if $B = (0)$.

- (b) Show that if B is also positive definite and $A - B$ is nonnegative definite, then $|A| \geq |B|$ with equality if and only if $A = B$.
- 3.53** Suppose that A is an $m \times m$ symmetric matrix with eigenvalues $\lambda_1, \dots, \lambda_m$ and associated eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$, whereas B is an $m \times m$ symmetric matrix with eigenvalues $\gamma_1, \dots, \gamma_m$ and associated eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$; that is, A and B have common eigenvectors.
- (a) Find the eigenvalues and eigenvectors of $A + B$.
 - (b) Find the eigenvalues and eigenvectors of AB .
 - (c) Show that $AB = BA$.
- 3.54** Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_r$ is a set of orthonormal eigenvectors corresponding to the r largest eigenvalues $\gamma_1, \dots, \gamma_r$ of the $m \times m$ symmetric matrix A and assume that $\gamma_r > \gamma_{r+1}$. Let P be the total eigenprojection of A associated with the eigenvalues $\gamma_1, \dots, \gamma_r$; that is,

$$P = \sum_{i=1}^r \mathbf{x}_i \mathbf{x}_i'.$$

Let B be another $m \times m$ symmetric matrix with its r largest eigenvalues given by μ_1, \dots, μ_r , where $\mu_r > \mu_{r+1}$, and a corresponding set of orthonormal eigenvectors given by $\mathbf{y}_1, \dots, \mathbf{y}_r$. Let Q be the total eigenprojection of B associated with the eigenvalues μ_1, \dots, μ_r so that

$$Q = \sum_{i=1}^r \mathbf{y}_i \mathbf{y}_i'.$$

(a) Show that $P = Q$ if and only if

$$\sum_{i=1}^r \{\gamma_i + \mu_i - \lambda_i(A + B)\} = 0.$$

(b) Let $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, where $\mathbf{x}_{r+1}, \dots, \mathbf{x}_m$ is a set of orthonormal eigenvectors corresponding to the smallest $m - r$ eigenvalues of A . Show that if $P = Q$, then $X'BX$ has the block diagonal form

$$\begin{bmatrix} U & (0) \\ (0) & V \end{bmatrix},$$

where U is $r \times r$ and V is $(m - r) \times (m - r)$.

3.55 Let $\lambda_1 \geq \dots \geq \lambda_m$ be the eigenvalues of the $m \times m$ symmetric matrix A and $\mathbf{x}_1, \dots, \mathbf{x}_m$ be a set of corresponding orthonormal eigenvectors. For some k , define the total eigenprojection associated with the eigenvalues $\lambda_k, \dots, \lambda_m$ as

$$P = \sum_{i=k}^m \mathbf{x}_i \mathbf{x}_i'.$$

Show that $\lambda_k = \dots = \lambda_m = \lambda$ if and only if

$$P(A - \lambda I_m)P = (0).$$

3.56 Let A_1, \dots, A_k be $m \times m$ symmetric matrices, and let τ_i be one of the eigenvalues of A_i . Let $\mathbf{x}_1, \dots, \mathbf{x}_r$ be a set of orthonormal $m \times 1$ vectors, and define

$$P = \sum_{i=1}^r \mathbf{x}_i \mathbf{x}_i'.$$

Show that if each of the eigenvalues τ_i has multiplicity r and has $\mathbf{x}_1, \dots, \mathbf{x}_r$ as associated eigenvectors, then

$$P \left\{ \sum_{i=1}^k (A_i - \tau_i I_m)^2 \right\} P = (0).$$

3.57 Let $\lambda_1 \geq \dots \geq \lambda_m$ be the eigenvalues of the $m \times m$ symmetric matrix A .

(a) If B is an $m \times r$ matrix, show that

$$\min_{B'B=I_r} \text{tr}(B'AB) = \sum_{i=1}^r \lambda_{m-i+1}$$

and

$$\max_{B'B=I_r} \text{tr}(B'AB) = \sum_{i=1}^r \lambda_i.$$

(b) For $r = 1, \dots, m$, show that

$$\sum_{i=1}^r \lambda_{m-i+1} \leq \sum_{i=1}^r a_{ii} \leq \sum_{i=1}^r \lambda_i.$$

3.58 Let $\lambda_1 \geq \dots \geq \lambda_m$ be the eigenvalues of the $m \times m$ positive definite matrix A .

(a) If B is an $m \times r$ matrix, show that

$$\min_{B'B=I_r} |B'AB| = \prod_{i=1}^r \lambda_{m-i+1}$$

and

$$\max_{B'B=I_r} |B'AB| = \prod_{i=1}^r \lambda_i.$$

(b) Let A_r be the $r \times r$ submatrix of A consisting of the first r rows and first r columns of A . For $r = 1, \dots, m$, show that

$$\prod_{i=1}^r \lambda_{m-i+1} \leq |A_r| \leq \prod_{i=1}^r \lambda_i.$$

3.59 Let A be an $m \times m$ nonnegative definite matrix, whereas B and C are $m \times m$ positive definite matrices. If i, j , and k are integers between 1 and m inclusive and satisfying $j + k \leq i + 1$, show that

$$(a) \quad \lambda_i(AB) \leq \lambda_j(AC^{-1})\lambda_k(CB),$$

$$(b) \quad \lambda_{m-i+1}(AB) \geq \lambda_{m-j+1}(AC^{-1})\lambda_{m-k+1}(CB).$$

3.60 Let A and B be $m \times m$ positive definite matrices. Show that

$$\frac{\lambda_i^2(AB)}{\lambda_1(A)\lambda_1(B)} \leq \lambda_i(A)\lambda_i(B) \leq \frac{\lambda_i^2(AB)}{\lambda_m(A)\lambda_m(B)},$$

for $i = 1, \dots, m$.

3.61 Let A be an $m \times m$ positive definite matrix with eigenvalues $\lambda_1 \geq \dots \geq \lambda_m$. If X is an $m \times k$ matrix with $k \leq m/2$, show that

$$\min_{X'X=I_k} \frac{|X'AX|}{\sqrt{|X'A^2X|}} = \prod_{i=1}^k \frac{2\sqrt{\lambda_i\lambda_{m-i+1}}}{\lambda_i + \lambda_{m-i+1}}.$$

3.62 Consider the function $\psi(\mathbf{x})$ given in Section 3.9, where the $m \times m$ matrix A is positive definite.

- (a) Show that the stationary points of $\psi(\mathbf{x})$ subject to the constraint $\mathbf{x}'\mathbf{x} = 1$ satisfy

$$\frac{A^2\mathbf{x}}{\mathbf{x}'A^2\mathbf{x}} - \frac{2A\mathbf{x}}{\mathbf{x}'A\mathbf{x}} + \mathbf{x} = \mathbf{0}.$$

- (b) Show that the equation given in part (a) holds for all normalized eigenvectors of A and all normalized antieigenvectors of A .

4

MATRIX FACTORIZATIONS AND MATRIX NORMS

4.1 INTRODUCTION

In this chapter, we take a look at some useful ways of expressing a given matrix A in the form of a product of other matrices having some special structure or canonical form. In many applications, such a decomposition of A may reveal the key features of A that are of interest to us. These factorizations are particularly useful in multivariate distribution theory in that they can expedite the mathematical development and often simplify the generalization of results from a special case to a more general situation. Our focus here will be on conditions for the existence of these factorizations as well as mathematical properties and consequences of the factorizations. Details on the numerical computation of the component matrices in these factorizations can be found in texts on numerical methods. Some useful references are Golub and Van Loan (2013), Press, et al. (2007), and Stewart (1998, 2001).

4.2 THE SINGULAR VALUE DECOMPOSITION

The first factorization that we consider, the singular value decomposition, could be described as the most useful because this factorization is for a matrix of any size; the subsequent decompositions will only apply to square matrices. We will find this decomposition particularly useful in the next chapter when we generalize the concept of an inverse of a nonsingular square matrix to any matrix.

Theorem 4.1 If A is an $m \times n$ matrix of rank $r > 0$, orthogonal $m \times m$ and $n \times n$ matrices P and Q exist, such that $A = PDQ'$ and $D = P'AQ$, where the $m \times n$ matrix D is given by

$$\begin{aligned} \text{(a)} \quad \Delta \quad & \text{if} \quad r = m = n, & \text{(b)} \quad \begin{bmatrix} \Delta & (0) \end{bmatrix} \quad & \text{if} \quad r = m < n, \\ \text{(c)} \quad \begin{bmatrix} \Delta \\ (0) \end{bmatrix} \quad & \text{if} \quad r = n < m, & \text{(d)} \quad \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} \quad & \text{if} \quad r < m, r < n, \end{aligned}$$

and Δ is an $r \times r$ diagonal matrix with positive diagonal elements. The diagonal elements of Δ^2 are the positive eigenvalues of $A'A$ and AA' .

Proof. We will prove the result for the case $r < m$ and $r < n$. The proofs of (a)–(c) only require notational changes. Let Δ^2 be the $r \times r$ diagonal matrix whose diagonal elements are the r positive eigenvalues of $A'A$, which are identical to the positive eigenvalues of AA' by Theorem 3.27. Define Δ to be the diagonal matrix whose diagonal elements are the positive square roots of the corresponding diagonal elements of Δ^2 . Since $A'A$ is an $n \times n$ symmetric matrix, we can find an $n \times n$ orthogonal matrix Q , such that

$$Q'A'AQ = \begin{bmatrix} \Delta^2 & (0) \\ (0) & (0) \end{bmatrix}.$$

Partitioning Q as $Q = [Q_1 \quad Q_2]$, where Q_1 is $n \times r$, the identity above implies that

$$Q_1'A'AQ_1 = \Delta^2 \tag{4.1}$$

and

$$Q_2'A'AQ_2 = (0). \tag{4.2}$$

Note that from (4.2), it follows that

$$AQ_2 = (0). \tag{4.3}$$

Now let $P = [P_1 \quad P_2]$ be an $m \times m$ orthogonal matrix, where the $m \times r$ matrix $P_1 = AQ_1\Delta^{-1}$ and the $m \times (m-r)$ matrix P_2 is any matrix that makes P orthogonal. Consequently, we must have $P_2'P_1 = P_2'AQ_1\Delta^{-1} = (0)$ or, equivalently,

$$P_2'AQ_1 = (0). \tag{4.4}$$

By using (4.1), (4.3), and (4.4), we find that

$$P'AQ = \begin{bmatrix} P_1'AQ_1 & P_1'AQ_2 \\ P_2'AQ_1 & P_2'AQ_2 \end{bmatrix}$$

$$\begin{aligned}
&= \begin{bmatrix} \Delta^{-1}Q'_1A'AQ_1 & \Delta^{-1}Q'_1A'AQ_2 \\ P'_2AQ_1 & P'_2AQ_2 \end{bmatrix} \\
&= \begin{bmatrix} \Delta^{-1}\Delta^2 & \Delta^{-1}Q'_1A'(0) \\ (0) & P'_2(0) \end{bmatrix} \\
&= \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix},
\end{aligned}$$

and so the proof is complete. \square

The singular value decomposition as given in Theorem 4.1 is for real matrices A . However, this decomposition easily extends to complex matrices. In particular, an $m \times n$ complex matrix A can be expressed as $A = PDQ^*$, where the $m \times m$ and $n \times n$ matrices P and Q are unitary, while D has the same form given in Theorem 4.1 with the diagonal elements of Δ given by the positive square roots of the nonzero eigenvalues of A^*A .

The diagonal elements of the Δ matrix given in Theorem 4.1, that is, the positive square roots of the positive eigenvalues of $A'A$ and AA' , are called the singular values of A . It is obvious from the proof of Theorem 4.1 that the columns of Q form an orthonormal set of eigenvectors of $A'A$, and so

$$A'A = QD'DQ'. \quad (4.5)$$

It is important to note also that the columns of P form an orthonormal set of eigenvectors of AA' because

$$AA' = PDQ'QD'P' = PDD'P'. \quad (4.6)$$

If we again partition P and Q as $P = [P_1 \ P_2]$ and $Q = [Q_1 \ Q_2]$, where P_1 is $m \times r$ and Q_1 is $n \times r$, then the singular value decomposition can be restated as follows.

Corollary 4.1.1 If A is an $m \times n$ matrix of rank $r > 0$, then $m \times r$ and $n \times r$ matrices P_1 and Q_1 exist, such that $P'_1P_1 = Q'_1Q_1 = I_r$ and $A = P_1\Delta Q'_1$, where Δ is an $r \times r$ diagonal matrix with positive diagonal elements.

It follows from (4.5) and (4.6) that P_1 and Q_1 are semiorthogonal matrices satisfying

$$P'_1AA'P_1 = \Delta^2, \quad Q'_1A'AQ_1 = \Delta^2. \quad (4.7)$$

However, in the decomposition $A = P_1\Delta Q'_1$, the choice of the semiorthogonal matrix P_1 satisfying (4.7) is dependent on the choice of the Q_1 matrix. This should be apparent from the proof of Theorem 4.1 in which any semiorthogonal matrix Q_1 satisfying

(4.7) was first selected, but the choice of P_1 was then given by $P_1 = AQ_1\Delta^{-1}$. Alternatively, we could have first selected a semiorthogonal matrix P_1 satisfying (4.7) and then have chosen $Q_1 = A'P_1\Delta^{-1}$.

A lot of information about the structure of a matrix A can be obtained from its singular value decomposition. The number of singular values gives the rank of A , whereas the columns of P_1 and Q_1 are orthonormal bases for the column space and row space of A , respectively. Similarly, the columns of P_2 span the null space of A' , and the columns of Q_2 span the null space of A .

Theorem 4.1 and Corollary 4.1.1 are related to Theorem 1.11 and its corollary, Corollary 1.11.1, which were stated as consequences of the properties of elementary transformations. It is easily verified that Theorem 1.11 and Corollary 1.11.1 also follow directly from Theorem 4.1 and Corollary 4.1.1.

Example 4.1 We will find the singular value decomposition for the 4×3 matrix

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

First, an eigenanalysis of the matrix

$$A'A = \begin{bmatrix} 18 & -10 & 4 \\ -10 & 18 & 4 \\ 4 & 4 & 4 \end{bmatrix}$$

reveals that it has eigenvalues 28, 12, and 0 with associated normalized eigenvectors $(1/\sqrt{2}, -1/\sqrt{2}, 0)'$, $(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})'$, and $(1/\sqrt{6}, 1/\sqrt{6}, -2/\sqrt{6})'$, respectively. Let these eigenvectors be the columns of the 3×3 orthogonal matrix Q . Clearly, $\text{rank}(A) = 2$ and the two singular values of A are $\sqrt{28}$ and $\sqrt{12}$. Thus, the 4×2 matrix P_1 is given by

$$\begin{aligned} P_1 &= AQ_1\Delta^{-1} = \begin{bmatrix} 2 & 0 & 1 \\ 3 & -1 & 1 \\ -2 & 4 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{3} \\ -1/\sqrt{2} & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{bmatrix} \\ &\quad \times \begin{bmatrix} 1/\sqrt{28} & 0 \\ 0 & 1/\sqrt{12} \end{bmatrix} \\ &= \begin{bmatrix} 1/\sqrt{14} & 1/2 \\ 2/\sqrt{14} & 1/2 \\ -3/\sqrt{14} & 1/2 \\ 0 & 1/2 \end{bmatrix}. \end{aligned}$$

The 4×2 matrix P_2 can be any matrix satisfying $P_1'P_2 = (0)$ and $P_2'P_2 = I_2$; for instance, $(1/\sqrt{12}, 1/\sqrt{12}, 1/\sqrt{12}, -3/\sqrt{12})'$ and $(-5/\sqrt{42}, 4/\sqrt{42}, 1/\sqrt{42}, 0)'$

can be chosen as the columns of P_2 . Then our singular value decomposition of A is given by

$$\begin{bmatrix} 1/\sqrt{14} & 1/2 & 1/\sqrt{12} & -5/\sqrt{42} \\ 2/\sqrt{14} & 1/2 & 1/\sqrt{12} & 4/\sqrt{42} \\ -3/\sqrt{14} & 1/2 & 1/\sqrt{12} & 1/\sqrt{42} \\ 0 & 1/2 & -3/\sqrt{12} & 0 \end{bmatrix} \begin{bmatrix} \sqrt{28} & 0 & 0 \\ 0 & \sqrt{12} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \times \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{6} & 1/\sqrt{6} & -2/\sqrt{6} \end{bmatrix},$$

or in the form of Corollary 4.1.1,

$$\begin{bmatrix} 1/\sqrt{14} & 1/2 \\ 2/\sqrt{14} & 1/2 \\ -3/\sqrt{14} & 1/2 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} \sqrt{28} & 0 \\ 0 & \sqrt{12} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix}.$$

As an alternative way of determining the matrix P , we could have used the fact that its columns are eigenvectors of the matrix

$$AA' = \begin{bmatrix} 5 & 7 & -3 & 3 \\ 7 & 11 & -9 & 3 \\ -3 & -9 & 21 & 3 \\ 3 & 3 & 3 & 3 \end{bmatrix}.$$

However, when constructing P this way, one must check the decomposition $A = P_1 \Delta Q_1'$ to determine the correct sign for each of the columns of P_1 .

The singular value decomposition of a vector is very easy to construct. We illustrate this in the next example.

Example 4.2 Let \mathbf{x} be an $m \times 1$ nonnull vector. Its singular value decomposition will be of the form

$$\mathbf{x} = P\mathbf{d}q,$$

where P is an $m \times m$ orthogonal matrix, \mathbf{d} is an $m \times 1$ vector having only its first component nonzero, and q is a scalar satisfying $q^2 = 1$. The single singular value of \mathbf{x} is given by λ , where $\lambda^2 = \mathbf{x}'\mathbf{x}$. If we define $\mathbf{x}_* = \lambda^{-1}\mathbf{x}$, note that $\mathbf{x}_*'\mathbf{x}_* = 1$, and

$$\mathbf{x}\mathbf{x}'\mathbf{x}_* = \mathbf{x}\mathbf{x}'(\lambda^{-1}\mathbf{x}) = (\lambda^{-1}\mathbf{x})\mathbf{x}'\mathbf{x} = \lambda^2\mathbf{x}_*,$$

so that \mathbf{x}_* is a normalized eigenvector of $\mathbf{x}\mathbf{x}'$ corresponding to its single positive eigenvalue λ^2 . Any nonnull vector orthogonal to \mathbf{x}_* is an eigenvector of $\mathbf{x}\mathbf{x}'$

corresponding to the repeated eigenvalue 0. Thus, if we let $\mathbf{d} = (\lambda, 0, \dots, 0)'$, $q = 1$, and $P = (\mathbf{x}_*, \mathbf{p}_2, \dots, \mathbf{p}_m)$ be any orthogonal matrix with \mathbf{x}_* as its first column, then

$$P\mathbf{d}q = [\mathbf{x}_*, \mathbf{p}_2, \dots, \mathbf{p}_m] \begin{bmatrix} \lambda \\ 0 \\ \vdots \\ 0 \end{bmatrix} = \lambda \mathbf{x}_* = \mathbf{x},$$

as is required.

When A is $m \times m$ and symmetric, the singular values of A are directly related to the eigenvalues of A . This follows from the fact that $AA' = A^2$, and the eigenvalues of A^2 are the squares of the eigenvalues of A . Thus, the singular values of A will be given by the absolute values of the eigenvalues of A . If we let the columns of P be a set of orthonormal eigenvectors of A , then the Q matrix in Theorem 4.1 will be identical to P except that any column of Q that is associated with a negative eigenvalue will be -1 times the corresponding column of P . If A is nonnegative definite, then the singular values of A will be the same as the positive eigenvalues of A and, in fact, the singular value decomposition of A is simply the spectral decomposition of A discussed in the next section. This nice relationship between the eigenvalues and singular values of a symmetric matrix does not carry over to general square matrices.

Example 4.3 Consider the 2×2 matrix

$$A = \begin{bmatrix} 6 & 6 \\ -1 & 1 \end{bmatrix},$$

which has

$$AA' = \begin{bmatrix} 72 & 0 \\ 0 & 2 \end{bmatrix}, \quad A'A = \begin{bmatrix} 37 & 35 \\ 35 & 37 \end{bmatrix}.$$

Clearly, the singular values of A are $\sqrt{72} = 6\sqrt{2}$ and $\sqrt{2}$. Normalized eigenvectors corresponding to 72 and 2 are $(1, 0)'$ and $(0, 1)'$ for AA' , whereas $A'A$ has $(1/\sqrt{2}, 1/\sqrt{2})'$ and $(-1/\sqrt{2}, 1/\sqrt{2})'$. Thus, the singular value decomposition of A can be written as

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} 6\sqrt{2} & 0 \\ 0 & \sqrt{2} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{2} \\ -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}.$$

On the other hand, an eigenanalysis of A yields the eigenvalues 4 and 3. Associated normalized eigenvectors are $(3/\sqrt{10}, -1/\sqrt{10})'$ and $(2/\sqrt{5}, -1/\sqrt{5})'$.

We end this section with an example that illustrates an application of the singular value decomposition to least squares regression. For more discussion of this and other applications of the singular value decomposition in statistics, see Mandel (1982), Eubank and Webster (1985), and Nelder (1985).

Example 4.4 In this example, we will take a closer look at the multicollinearity problem that we first discussed in Example 3.7. Suppose that we have the standardized regression model,

$$\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}.$$

We have seen in Example 2.16 that the least squares estimator of δ_0 is \bar{y} . The fitted model $\hat{\mathbf{y}} = \bar{y} \mathbf{1}_N + Z_1 \hat{\boldsymbol{\delta}}_1$ gives points on a hyperplane in R^{k+1} , where the $(k+1)$ axes correspond to the k standardized explanatory variables and the fitted response variable. Now let $Z_1 = VDU'$ be the singular value decomposition of the $N \times k$ matrix Z_1 . Thus, V is an $N \times N$ orthogonal matrix, U is a $k \times k$ orthogonal matrix, and D is an $N \times k$ matrix that has the square roots of the eigenvalues of $Z_1' Z_1$ as its diagonal elements and zeros elsewhere. We can rewrite the model $\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}$ as we did in Example 2.16 by defining $\alpha_0 = \delta_0$, $\boldsymbol{\alpha}_1 = U' \boldsymbol{\delta}_1$ and $W_1 = VD$, so that $\mathbf{y} = \alpha_0 \mathbf{1}_N + W_1 \boldsymbol{\alpha}_1 + \boldsymbol{\epsilon}$. Suppose that exactly r of the diagonal elements of D , specifically the last r diagonal elements, are zeros, and so by partitioning U , V , and D appropriately, we get $Z_1 = V_1 D_1 U_1'$, where D_1 is a $(k-r) \times (k-r)$ diagonal matrix. This means that the row space of Z_1 is a $(k-r)$ -dimensional subspace of R^k , and this subspace is spanned by the columns of U_1 ; that is, the points on the fitted regression hyperplane described above, when projected onto the k -dimensional standardized explanatory variable space, are actually confined to a $(k-r)$ -dimensional subspace. Also, the model $\mathbf{y} = \alpha_0 \mathbf{1}_N + W_1 \boldsymbol{\alpha}_1 + \boldsymbol{\epsilon}$ simplifies to

$$\mathbf{y} = \alpha_0 \mathbf{1}_N + W_{11} \boldsymbol{\alpha}_{11} + \boldsymbol{\epsilon}, \quad (4.8)$$

where $W_{11} = V_1 D_1$, $\boldsymbol{\alpha}_{11} = U_1' \boldsymbol{\delta}_1$, and the least squares estimator of the $(k-r) \times 1$ vector $\boldsymbol{\alpha}_{11}$ is given by $\hat{\boldsymbol{\alpha}}_{11} = (W_{11}' W_{11})^{-1} W_{11}' \mathbf{y} = D_1^{-1} V_1' \mathbf{y}$. This can be used to find a least squares estimator of $\boldsymbol{\delta}_1$ because we must have $\hat{\boldsymbol{\alpha}}_{11} = U_1' \hat{\boldsymbol{\delta}}_1$. Partitioning $\hat{\boldsymbol{\delta}}_1 = (\hat{\boldsymbol{\delta}}_{11}', \hat{\boldsymbol{\delta}}_{12}')'$ and $U_1' = (U_{11}', U_{12}')$, where $\hat{\boldsymbol{\delta}}_{11}$ is $(k-r) \times 1$, we obtain the relationship

$$\hat{\boldsymbol{\alpha}}_{11} = U_{11}' \hat{\boldsymbol{\delta}}_{11} + U_{12}' \hat{\boldsymbol{\delta}}_{12}.$$

Premultiplying this equation by $U_{11}^{\prime -1}$ (if U_{11} is not nonsingular, then $\boldsymbol{\delta}_1$ and U_1 can be rearranged so that it is), we find that

$$\hat{\boldsymbol{\delta}}_{11} = U_{11}^{\prime -1} \hat{\boldsymbol{\alpha}}_{11} - U_{11}^{\prime -1} U_{12}' \hat{\boldsymbol{\delta}}_{12};$$

that is, the least squares estimator of $\boldsymbol{\delta}_1$ is not unique because $\hat{\boldsymbol{\delta}}_1 = (\hat{\boldsymbol{\delta}}_{11}', \hat{\boldsymbol{\delta}}_{12}')'$ is a least squares estimator for any choice of $\hat{\boldsymbol{\delta}}_{12}$, as long as $\hat{\boldsymbol{\delta}}_{11}$ satisfies this identity. Now suppose that we wish to estimate the response variable y corresponding to an observation that has the standardized explanatory variables at the values given in the $k \times 1$ vector \mathbf{z} . Using a least squares estimate $\hat{\boldsymbol{\delta}}_1$ we obtain the estimate $\hat{y} = \bar{y} + \mathbf{z}' \hat{\boldsymbol{\delta}}_1$. This estimated response, like $\hat{\boldsymbol{\delta}}_1$, may not be unique because, if we partition \mathbf{z} as $\mathbf{z}' = (\mathbf{z}_1', \mathbf{z}_2')'$ with \mathbf{z}_1 being $(k-r) \times 1$,

$$\begin{aligned} \hat{y} &= \bar{y} + \mathbf{z}' \hat{\boldsymbol{\delta}}_1 = \bar{y} + \mathbf{z}_1' \hat{\boldsymbol{\delta}}_{11} + \mathbf{z}_2' \hat{\boldsymbol{\delta}}_{12} \\ &= \bar{y} + \mathbf{z}_1' U_{11}^{\prime -1} \hat{\boldsymbol{\alpha}}_{11} + (\mathbf{z}_2' - \mathbf{z}_1' U_{11}^{\prime -1} U_{12}') \hat{\boldsymbol{\delta}}_{12}. \end{aligned}$$

Thus, \hat{y} does not depend on the arbitrary $\hat{\delta}_{12}$ and is therefore unique, only if

$$(z'_2 - z'_1 U'^{-1}_{11} U'_{12}) = \mathbf{0}', \quad (4.9)$$

in which case the unique estimated value is given by $\hat{y} = \bar{y} + z'_1 U'^{-1}_{11} \hat{\alpha}_{11}$. It is easily shown that the set of all vectors $z = (z'_1, z'_2)'$ satisfying (4.9) is simply the column space of U_1 . Thus, $y = \delta_0 + z' \delta_1$ is uniquely estimated only if the vector of standardized explanatory variables z falls within the space spanned by the collection of all vectors of standardized explanatory variables available to compute $\hat{\delta}_1$.

In the typical multicollinearity problem, Z_1 is full rank so that the matrix D has no zero diagonal elements but instead has r of its diagonal elements very small relative to the others. In this case, the row space of Z_1 is all of R^k , but the points corresponding to the rows of Z_1 all lie very close to a $(k - r)$ -dimensional subspace S of R^k , specifically, the space spanned by the columns of U_1 . Small changes in the values of the response variables corresponding to these points can substantially alter the position of the fitted regression hyperplane $\hat{y} = \bar{y} + z' \hat{\delta}_1$ for vectors z lying outside of and, in particular, far from S . For instance, if $k = 2$ and $r = 1$, the points corresponding to the rows of Z_1 all lie very close to S , which, in this case, is a line in the z_1, z_2 plane, and $\hat{y} = \bar{y} + z' \hat{\delta}_1$ will be given by a plane in R^3 extended over the z_1, z_2 plane. The fitted regression plane $\hat{y} = \bar{y} + z' \hat{\delta}_1$ can be identified by the line formed as the intersection of this plane and the plane perpendicular to the z_1, z_2 plane and passing through the line S , along with the tilt of the fitted regression plane. Small changes in the values of the response variables will produce small changes in both the location of this line of intersection and the tilt of the plane. However, even a slight change in the tilt of the regression plane will yield large changes on the surface of this plane for vectors z far from S . The adverse effect of this tilting can be eliminated by the use of principal components regression. As we saw in Example 3.7, principal components regression utilizes the regression model (4.8), and so an estimated response will be given by $\hat{y} = \bar{y} + z' U_1 D_1^{-1} V_1' y$. Since this regression model technically holds only for $z \in S$, by using this model for $z \notin S$, we will introduce bias into our estimate of y . The advantage of principal components regression is that this bias may be compensated for by a large enough reduction in the variance of our estimate so as to reduce the mean squared error (see Problem 4.11). However, it should be apparent that the predicted values of y obtained from both ordinary least squares regression and principal components regression will be poor if the vector z is far from S .

4.3 THE SPECTRAL DECOMPOSITION OF A SYMMETRIC MATRIX

The spectral decomposition of a symmetric matrix, which was briefly discussed in Chapter 3, is nothing more than a special case of the singular value decomposition. We summarize this result in Theorem 4.2.

Theorem 4.2 Let A be an $m \times m$ symmetric matrix with eigenvalues $\lambda_1, \dots, \lambda_m$, and suppose that x_1, \dots, x_m is a set of orthonormal eigenvectors corresponding

to these eigenvalues. Then, if $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ and $X = (x_1, \dots, x_m)$, it follows that

$$A = X\Lambda X'.$$

We can use the spectral decomposition of a nonnegative definite matrix A to find a square root matrix of A ; that is, we wish to find an $m \times m$ nonnegative definite matrix $A^{1/2}$ for which $A = A^{1/2}A^{1/2}$. If Λ and X are defined as in Theorem 4.2, and we let $\Lambda^{1/2} = \text{diag}(\lambda_1^{1/2}, \dots, \lambda_m^{1/2})$ and $A^{1/2} = X\Lambda^{1/2}X'$, then because $X'X = I_m$,

$$\begin{aligned} A^{1/2}A^{1/2} &= X\Lambda^{1/2}X'X\Lambda^{1/2}X' = X\Lambda^{1/2}\Lambda^{1/2}X' \\ &= X\Lambda X' = A, \end{aligned}$$

as is required. Note that $(A^{1/2})' = (X\Lambda^{1/2}X')' = X\Lambda^{1/2}X' = A^{1/2}$; consequently, $X\Lambda^{1/2}X'$ is referred to as the symmetric square root of A . Note also that if we did not require A to be nonnegative definite, then $A^{1/2}$ would be a complex matrix if some of the eigenvalues of A are negative.

It is easy to show that the nonnegative definite square root $A^{1/2}$ is uniquely defined. In Section 3.4, we saw that the spectral decomposition of A can be written as

$$A = \sum_{i=1}^k \mu_i P_A(\mu_i),$$

where μ_1, \dots, μ_k are the spectral values of A . The eigenprojections of A are uniquely defined and satisfy $\{P_A(\mu_i)\}' = P_A(\mu_i)$, $\{P_A(\mu_i)\}^2 = P_A(\mu_i)$, and $P_A(\mu_i)P_A(\mu_j) = (0)$ if $i \neq j$. Suppose B is another $m \times m$ nonnegative definite matrix with its spectral decomposition given by

$$B = \sum_{j=1}^r \gamma_j P_B(\gamma_j).$$

If $A = B^2$ so that

$$A = \sum_{j=1}^r \gamma_j^2 P_B(\gamma_j),$$

it follows that we must have $r = k$ and for each i , $\mu_i = \gamma_j^2$ and $P_A(\mu_i) = P_B(\gamma_j)$ for some j . This implies that

$$B = \sum_{i=1}^k \mu_i^{1/2} P_A(\mu_i),$$

and this is equivalent to the expression $A^{1/2} = X\Lambda^{1/2}X'$ given above.

We can expand the set of square root matrices if we do not insist that $A^{1/2}$ be symmetric; that is, now let us consider any matrix $A^{1/2}$ satisfying $A = A^{1/2}(A^{1/2})'$. If Q is any $m \times m$ orthogonal matrix, then $A^{1/2} = X\Lambda^{1/2}Q'$ is such a square root

matrix because

$$\begin{aligned} A^{1/2} A^{1/2'} &= X \Lambda^{1/2} Q' Q \Lambda^{1/2} X' = X \Lambda^{1/2} \Lambda^{1/2} X' \\ &= X \Lambda X' = A. \end{aligned}$$

If $A^{1/2}$ is a lower triangular matrix with nonnegative diagonal elements, then the factorization $A = A^{1/2} A^{1/2'}$ is known as the Cholesky decomposition of A . Theorem 4.3 establishes the existence of such a decomposition.

Theorem 4.3 Let A be an $m \times m$ nonnegative definite matrix. Then an $m \times m$ lower triangular matrix T having nonnegative diagonal elements exists, such that $A = TT'$. Further, if A is positive definite, then the matrix T is unique and has positive diagonal elements.

Proof. We will prove the result for positive definite matrices. Our proof is by induction. The result clearly holds if $m = 1$, because in this case, A is a positive scalar, and so the unique T would be given by the positive square root of A . Now assume that the result holds for all positive definite $(m - 1) \times (m - 1)$ matrices. Partition A as

$$A = \begin{bmatrix} A_{11} & \mathbf{a}_{12} \\ \mathbf{a}'_{12} & a_{22} \end{bmatrix},$$

where A_{11} is $(m - 1) \times (m - 1)$. Since A_{11} must be positive definite if A is, we know there exists a unique $(m - 1) \times (m - 1)$ lower triangular matrix T_{11} having positive diagonal elements and satisfying $A_{11} = T_{11} T'_{11}$. Our proof will be complete if we can show that there is a unique $(m - 1) \times 1$ vector \mathbf{t}_{12} and a unique positive scalar t_{22} , such that

$$\begin{aligned} \begin{bmatrix} A_{11} & \mathbf{a}_{12} \\ \mathbf{a}'_{12} & a_{22} \end{bmatrix} &= \begin{bmatrix} T_{11} & \mathbf{0} \\ \mathbf{t}'_{12} & t_{22} \end{bmatrix} \begin{bmatrix} T'_{11} & \mathbf{t}_{12} \\ \mathbf{0}' & t_{22} \end{bmatrix} \\ &= \begin{bmatrix} T_{11} T'_{11} & T_{11} \mathbf{t}_{12} \\ \mathbf{t}'_{12} T'_{11} & \mathbf{t}'_{12} \mathbf{t}_{12} + t_{22}^2 \end{bmatrix}; \end{aligned}$$

that is, we must have $\mathbf{a}_{12} = T_{11} \mathbf{t}_{12}$ and $a_{22} = \mathbf{t}'_{12} \mathbf{t}_{12} + t_{22}^2$. Since T_{11} must be nonsingular, the unique choice of \mathbf{t}_{12} is given by $\mathbf{t}_{12} = T_{11}^{-1} \mathbf{a}_{12}$, and so t_{22}^2 must satisfy

$$\begin{aligned} t_{22}^2 &= a_{22} - \mathbf{t}'_{12} \mathbf{t}_{12} = a_{22} - \mathbf{a}'_{12} (T_{11}^{-1})' T_{11}^{-1} \mathbf{a}_{12} \\ &= a_{22} - \mathbf{a}'_{12} (T_{11} T'_{11})^{-1} \mathbf{a}_{12} = a_{22} - \mathbf{a}'_{12} A_{11}^{-1} \mathbf{a}_{12}. \end{aligned}$$

Note that because A is positive definite, $a_{22} - \mathbf{a}'_{12} A_{11}^{-1} \mathbf{a}_{12}$ will be positive because, if we let $\mathbf{x} = (\mathbf{x}'_1, -1)' = (\mathbf{a}'_{12} A_{11}^{-1}, -1)'$, then

$$\begin{aligned} \mathbf{x}' A \mathbf{x} &= \mathbf{x}'_1 A_{11} \mathbf{x}_1 - 2 \mathbf{x}'_1 \mathbf{a}_{12} + a_{22} \\ &= \mathbf{a}'_{12} A_{11}^{-1} A_{11} A_{11}^{-1} \mathbf{a}_{12} - 2 \mathbf{a}'_{12} A_{11}^{-1} \mathbf{a}_{12} + a_{22} \end{aligned}$$

$$= a_{22} - \mathbf{a}'_{12} A_{11}^{-1} \mathbf{a}_{12}.$$

Consequently, the unique $t_{22} > 0$ is given by $t_{22} = (a_{22} - \mathbf{a}'_{12} A_{11}^{-1} \mathbf{a}_{12})^{1/2}$. \square

The following decomposition, commonly known as the QR factorization, can be used to establish the triangular factorization of Theorem 4.3 for positive semidefinite matrices.

Theorem 4.4 Let A be an $m \times n$ matrix, where $m \geq n$. An $n \times n$ upper triangular matrix R with nonnegative diagonal elements and an $m \times n$ matrix Q satisfying $Q'Q = I_n$ exist, such that $A = QR$.

Proof. Let $r_{11} = (\mathbf{a}'_1 \mathbf{a}_1)^{1/2}$, where \mathbf{a}_1 is the first column of A . If $r_{11} = 0$, let $Q_1 = I_m$, otherwise let Q_1 be any $m \times m$ orthogonal matrix with its first row given by \mathbf{a}'_1/r_{11} so that

$$Q_1 A = \begin{bmatrix} r_{11} & \mathbf{b}'_{11} \\ \mathbf{0} & A_2 \end{bmatrix}$$

for some $(n-1) \times 1$ vector \mathbf{b}_{11} and $(m-1) \times (n-1)$ matrix A_2 . Let $r_{22} = (\mathbf{a}'_2 \mathbf{a}_2)^{1/2}$, where \mathbf{a}_2 is the first column of A_2 . If $r_{22} = 0$, let $Q_{22} = I_{m-1}$, otherwise let Q_{22} be any $(m-1) \times (m-1)$ orthogonal matrix with its first row given by \mathbf{a}'_2/r_{22} so that

$$Q_2 Q_1 A = \begin{bmatrix} r_{11} & r_{12} & \mathbf{b}'_{12} \\ 0 & r_{22} & \mathbf{b}'_{22} \\ \mathbf{0} & \mathbf{0} & A_3 \end{bmatrix}, \text{ where } Q_2 = \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & Q_{22} \end{bmatrix},$$

$\mathbf{b}'_{11} = (r_{12}, \mathbf{b}'_{12})$, \mathbf{b}_{22} is some $(n-2) \times 1$ vector, and A_3 is some $(m-2) \times (n-2)$ matrix. Continuing in this fashion through n steps, we will have obtained $m \times m$ orthogonal matrices Q_1, \dots, Q_n such that

$$Q_n \cdots Q_1 A = \begin{bmatrix} R \\ (0) \end{bmatrix},$$

where R is an $n \times n$ upper triangular matrix with diagonal elements r_{11}, \dots, r_{nn} , all of which are nonnegative. Finally, partition the orthogonal matrix $Q_* = Q_n \cdots Q_1$ as $Q_* = [Q, P]$, where Q is $m \times n$ so that $Q'Q = I_n$ and $A = QR$. \square

If A is a positive semidefinite matrix and $A = A^{1/2}(A^{1/2})'$, then the triangular factorization of Theorem 4.3 for positive semidefinite matrices can be proven by using the QR factorization of $(A^{1/2})'$.

Example 4.5 Suppose that the $m \times 1$ random vector \mathbf{x} has mean vector $\boldsymbol{\mu}$ and the positive definite covariance matrix Ω . By using a square root matrix of Ω , we can determine a linear transformation of \mathbf{x} so that the transformed random vector

is standardized; that is, it has mean vector $\mathbf{0}$ and covariance matrix I_m . If we let $\Omega^{1/2}$ be any matrix satisfying $\Omega = \Omega^{1/2}(\Omega^{1/2})'$ and put $\mathbf{z} = \Omega^{-1/2}(\mathbf{x} - \boldsymbol{\mu})$, where $\Omega^{-1/2} = (\Omega^{1/2})^{-1}$, then by using (1.8) and (1.9) of Section 1.13, we find that

$$\begin{aligned} E(\mathbf{z}) &= E\{\Omega^{-1/2}(\mathbf{x} - \boldsymbol{\mu})\} = \Omega^{-1/2}\{E(\mathbf{x} - \boldsymbol{\mu})\} \\ &= \Omega^{-1/2}(\boldsymbol{\mu} - \boldsymbol{\mu}) = \mathbf{0} \end{aligned}$$

and

$$\begin{aligned} \text{var}(\mathbf{z}) &= \text{var}\{\Omega^{-1/2}(\mathbf{x} - \boldsymbol{\mu})\} \\ &= \Omega^{-1/2}\{\text{var}(\mathbf{x} - \boldsymbol{\mu})\}(\Omega^{-1/2})' \\ &= \Omega^{-1/2}\{\text{var}(\mathbf{x})\}(\Omega^{-1/2})' \\ &= \Omega^{-1/2}\Omega(\Omega^{-1/2})' = I_m. \end{aligned}$$

Since the covariance matrix of \mathbf{z} is the identity matrix, the Euclidean distance function will give a meaningful measure of the distance between observations from this distribution. By using the linear transformation defined above, we can relate distances between \mathbf{z} observations to distances between \mathbf{x} observations. For example, the Euclidean distance between an observation \mathbf{z} and its expected value $\mathbf{0}$ is

$$\begin{aligned} d_I(\mathbf{z}, \mathbf{0}) &= \{(\mathbf{z} - \mathbf{0})'(\mathbf{z} - \mathbf{0})\}^{1/2} = (\mathbf{z}'\mathbf{z})^{1/2} \\ &= \{(\mathbf{x} - \boldsymbol{\mu})'(\Omega^{-1/2})'\Omega^{-1/2}(\mathbf{x} - \boldsymbol{\mu})\}^{1/2} \\ &= \{(\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\}^{1/2} \\ &= d_\Omega(\mathbf{x}, \boldsymbol{\mu}), \end{aligned}$$

where d_Ω is the Mahalanobis distance function defined in Section 2.2. Similarly, if \mathbf{x}_1 and \mathbf{x}_2 are two observations from the distribution of \mathbf{x} and \mathbf{z}_1 and \mathbf{z}_2 are the corresponding transformed vectors, then $d_I(\mathbf{z}_1, \mathbf{z}_2) = d_\Omega(\mathbf{x}_1, \mathbf{x}_2)$. This relationship between the Mahalanobis distance and the Euclidean distance makes the construction of the Mahalanobis distance function more apparent. The Mahalanobis distance is nothing more than a two-stage computation of distance; the first stage transforms points so as to remove the effect of correlations and differing variances, whereas the second stage simply computes the Euclidean distance for these transformed points.

Example 4.6 We consider the generalized least squares regression model which was first discussed in Example 2.22. This model has the same form as the standard multiple regression model that is,

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

but now $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 C$, where C is a known $N \times N$ positive definite matrix. In Example 2.22, we found the least squares estimator of $\boldsymbol{\beta}$ by using properties of projections. In this example, we obtain that same estimator by a different approach. We will transform the problem to ordinary least squares regression; that is, we wish to

transform the model so that the vector of random errors in the transformed model has $\sigma^2 I_N$ as its covariance matrix. This can be done by using any square root matrix of C . Let T be any $N \times N$ matrix satisfying $TT' = C$ or, equivalently, $T'^{-1}T^{-1} = C^{-1}$. Now transform our original regression model to the model

$$\mathbf{y}_* = X_*\beta + \epsilon_*,$$

where $\mathbf{y}_* = T^{-1}\mathbf{y}$, $X_* = T^{-1}X$, and $\epsilon_* = T^{-1}\epsilon$, and note that $E(\epsilon_*) = T^{-1}E(\epsilon) = \mathbf{0}$ and

$$\begin{aligned}\text{var}(\epsilon_*) &= \text{var}(T^{-1}\epsilon) = T^{-1}\{\text{var}(\epsilon)\}T'^{-1} \\ &= T^{-1}(\sigma^2 C)T'^{-1} = \sigma^2 T^{-1}TT'T'^{-1} \\ &= \sigma^2 I_N.\end{aligned}$$

Thus, the generalized least squares estimator $\hat{\beta}_*$ of β in the model $\mathbf{y} = X\beta + \epsilon$ is given by the ordinary least squares estimator of β in the model $\mathbf{y}_* = X_*\beta + \epsilon_*$, and so it can be expressed as

$$\begin{aligned}\hat{\beta}_* &= (X'_*X_*)^{-1}X'_*\mathbf{y}_* = (X'T'^{-1}T^{-1}X)^{-1}X'T'^{-1}T^{-1}\mathbf{y} \\ &= (X'C^{-1}X)^{-1}X'C^{-1}\mathbf{y}.\end{aligned}$$

In some situations, a matrix A can be expressed in the form of the transpose product, BB' , where the $m \times r$ matrix B has $r < m$, so that unlike a square root matrix, B is not square. This is the subject of Theorem 4.5, the proof of which will be left to the reader as an exercise.

Theorem 4.5 Let A be an $m \times m$ nonnegative definite matrix with $\text{rank}(A) = r$. Then there exists an $m \times r$ matrix B having rank of r , such that $A = BB'$.

The transpose product form $A = BB'$ of the nonnegative definite matrix A is not unique. However, if C is another matrix of order $m \times n$ where $n \geq r$ and $A = CC'$, then there is an explicit relationship between the matrices B and C . This is established in the next theorem.

Theorem 4.6 Suppose that B is an $m \times h$ matrix and C is an $m \times n$ matrix, where $h \leq n$. Then $BB' = CC'$ if and only if an $h \times n$ matrix Q exists, such that $QQ' = I_h$ and $C = BQ$.

Proof. If $C = BQ$ with $QQ' = I_h$, then clearly

$$CC' = BQ(BQ)' = BQQ'B' = BB'.$$

Conversely, now suppose that $BB' = CC'$. We will assume that $h = n$ because if $h < n$, we can form the matrix $B_* = [B \quad \mathbf{0}]$, so that B_* is $m \times n$ and

$B_*B'_* = BB'$; then proving that an $n \times n$ orthogonal matrix Q_* exists, such that $C = B_*Q_*$, will yield $C = BQ$ if we take Q to be the first h rows of Q_* . Now because BB' is symmetric, an orthogonal matrix X exists, such that

$$BB' = CC' = X \begin{bmatrix} \Lambda & (0) \\ (0) & (0) \end{bmatrix} X' = X_1 \Lambda X_1',$$

where $\text{rank}(BB') = r$ and the $r \times r$ diagonal matrix Λ contains the positive eigenvalues of the nonnegative definite matrix BB' . Here X has been partitioned as $X = [X_1 \ X_2]$, where X_1 is $m \times r$. Form the matrices

$$\begin{aligned} E &= \begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' B \\ &= \begin{bmatrix} \Lambda^{-1/2} X_1' B \\ X_2' B \end{bmatrix} = \begin{bmatrix} E_1 \\ E_2 \end{bmatrix}, \end{aligned} \quad (4.10)$$

$$\begin{aligned} F &= \begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' C \\ &= \begin{bmatrix} \Lambda^{-1/2} X_1' C \\ X_2' C \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}, \end{aligned} \quad (4.11)$$

so that

$$EE' = FF' = \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix};$$

that is, $E_1 E_1' = F_1 F_1' = I_r$, $E_2 E_2' = F_2 F_2' = (0)$, and so $E_2 = F_2 = (0)$. Now let E_3 and F_3 be any $(h-r) \times h$ matrices, such that $E_* = [E_1' \ E_3']'$ and $F_* = [F_1' \ F_3']'$ are both orthogonal matrices. Consequently, if $Q = E_*' F_*$, then $QQ' = E_*' F_*' F_* E_* = E_*' E_* = I_h$, so Q is orthogonal. Since E_* is orthogonal, we have $E_1 E_3' = (0)$, and so

$$\begin{aligned} EQ &= EE_*' F_* = \begin{bmatrix} E_1 \\ (0) \end{bmatrix} [E_1' \ E_3'] \begin{bmatrix} F_1 \\ F_3 \end{bmatrix} \\ &= \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} F_1 \\ F_3 \end{bmatrix} = \begin{bmatrix} F_1 \\ (0) \end{bmatrix} = F. \end{aligned}$$

However, using (4.10) and (4.11), $EQ = F$ can be written as

$$\begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' BQ = \begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' C.$$

The result now follows by premultiplying this equation by

$$X \begin{bmatrix} \Lambda^{1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix},$$

because $XX' = I_m$. □

4.4 THE DIAGONALIZATION OF A SQUARE MATRIX

From the spectral decomposition theorem, we know that every symmetric matrix can be transformed to a diagonal matrix by postmultiplying by an appropriately chosen orthogonal matrix and premultiplying by its transpose. This result gives us a very useful and simple relationship between a symmetric matrix and its eigenvalues and eigenvectors. In this section, we investigate a generalization of this relationship to square matrices in general. We begin with Definition 4.1.

Definition 4.1 The $m \times m$ matrices A and B are said to be similar matrices if a nonsingular matrix C exists, such that $A = CBC^{-1}$.

It follows from Theorem 3.2(d) that similar matrices have identical eigenvalues. However, the converse is not true. For instance, if we have

$$A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix},$$

then A and B have identical eigenvalues because each has 0 with multiplicity 2. Clearly, however, there is no nonsingular matrix C satisfying $A = CBC^{-1}$.

The spectral decomposition theorem given as Theorem 4.2 tells us that every symmetric matrix is similar to a diagonal matrix. Unfortunately, the same statement does not hold for all square matrices. If the diagonal elements of the diagonal matrix Λ are the eigenvalues of A , and the columns of X are corresponding eigenvectors, then the eigenvalue-eigenvector equation $AX = X\Lambda$ immediately leads to the identity $X^{-1}AX = \Lambda$, if X is nonsingular; that is, the diagonalization of an $m \times m$ matrix simply depends on the existence of a set of m linearly independent eigenvectors. Consequently, we have the following result, previously mentioned in Section 3.3, which follows immediately from Theorem 3.7.

Theorem 4.7 Suppose that the $m \times m$ matrix A has the eigenvalues $\lambda_1, \dots, \lambda_m$, which are distinct. If $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ and $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, where $\mathbf{x}_1, \dots, \mathbf{x}_m$ are eigenvectors of A corresponding to $\lambda_1, \dots, \lambda_m$, then

$$X^{-1}AX = \Lambda. \quad (4.12)$$

Theorem 4.7 gives a sufficient but not necessary condition for the diagonalization of a general square matrix; that is, some nonsymmetric matrices that have multiple eigenvalues are similar to a diagonal matrix. The next theorem gives a necessary and sufficient condition for a matrix to be diagonalizable.

Theorem 4.8 Suppose the eigenvalues $\lambda_1, \dots, \lambda_m$ of the $m \times m$ matrix A consist of h distinct values μ_1, \dots, μ_h having multiplicities r_1, \dots, r_h , so that $r_1 + \dots + r_h = m$. Then A has a set of m linearly independent eigenvectors and, thus, is diagonalizable if and only if $\text{rank}(A - \mu_i I_m) = m - r_i$ for $i = 1, \dots, h$.

Proof. First, suppose that A is diagonalizable, so that using the usual notation, we have $X^{-1}AX = \Lambda$ or, equivalently, $A = \Lambda X^{-1}$. Thus,

$$\begin{aligned} \text{rank}(A - \mu_i I_m) &= \text{rank}(X \Lambda X^{-1} - \mu_i I_m) \\ &= \text{rank}\{X(\Lambda - \mu_i I_m)X^{-1}\} \\ &= \text{rank}(\Lambda - \mu_i I_m), \end{aligned}$$

where the last equality follows from the fact that the rank of a matrix is unaltered by its multiplication by a nonsingular matrix. Now, because μ_i has multiplicity r_i , the diagonal matrix $(\Lambda - \mu_i I_m)$ has exactly $m - r_i$ nonzero diagonal elements, which then guarantees that $\text{rank}(A - \mu_i I_m) = m - r_i$. Conversely, now suppose that $\text{rank}(A - \mu_i I_m) = m - r_i$, for $i = 1, \dots, h$. This implies that the dimension of the null space of $(A - \mu_i I_m)$ is $m - (m - r_i) = r_i$, and so we can find r_i linearly independent vectors satisfying the equation

$$(A - \mu_i I_m)\mathbf{x} = \mathbf{0}.$$

However, any such \mathbf{x} is an eigenvector of A corresponding to the eigenvalue μ_i . Consequently, we can find a set of r_i linearly independent eigenvectors associated with the eigenvalue μ_i . From Theorem 3.7, we know that eigenvectors corresponding to different eigenvalues are linearly independent. As a result, any set of m eigenvectors of A , which has r_i linearly independent eigenvectors corresponding to μ_i for each i , will also be linearly independent. Therefore, A is diagonalizable, and so the proof is complete. \square

The spectral decomposition previously given for a symmetric matrix can be extended to diagonalizable matrices. Let A be an $m \times m$ diagonalizable matrix with eigenvalues $\lambda_1, \dots, \lambda_m$ and the corresponding linearly independent eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$. Denote the spectral set of A by $\{\mu_1, \dots, \mu_h\}$ with μ_i having multiplicity r_i and $\lambda_{M_i+1} = \dots = \lambda_{M_i+r_i} = \mu_i$, where $M_1 = 0$ and $M_i = \sum_{j=1}^{i-1} r_j$ for $i = 2, \dots, h$. If $Y = (\mathbf{y}_1, \dots, \mathbf{y}_m) = X^{-1'}$, then A can be expressed as

$$A = X \Lambda X^{-1} = X \Lambda Y' = \sum_{i=1}^m \lambda_i \mathbf{x}_i \mathbf{y}_i' = \sum_{i=1}^h \mu_i P_A(\mu_i),$$

where $P_A(\mu_i) = \sum_{j=1}^{r_i} \mathbf{x}_{M_i+j} \mathbf{y}_{M_i+j}'$. If A is symmetric, then X is an orthogonal matrix so $Y = X$ and this, of course, reduces to the spectral decomposition given in Section 4.3.

Example 4.7 In Example 3.3, we saw that the matrix

$$A = \begin{bmatrix} 2 & -1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

has eigenvalues $\lambda_1 = \lambda_2 = 1$ and $\lambda_3 = 2$, and corresponding linearly independent eigenvectors $\mathbf{x}_1 = (0, 0, 1)'$, $\mathbf{x}_2 = (1, 1, 0)'$, and $\mathbf{x}_3 = (1, 0, 0)'$. With $X = (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3)$, it is easily shown that

$$X^{-1} = \begin{bmatrix} \mathbf{y}'_1 \\ \mathbf{y}'_2 \\ \mathbf{y}'_3 \end{bmatrix} = \begin{bmatrix} 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & -1 & 0 \end{bmatrix}.$$

Thus, the spectral decomposition of A is

$$\begin{aligned} A &= 1P_A(1) + 2P_A(2) = 1(\mathbf{x}_1\mathbf{y}'_1 + \mathbf{x}_2\mathbf{y}'_2) + 2\mathbf{x}_3\mathbf{y}'_3 \\ &= 1 \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} + 2 \begin{bmatrix} 1 & -1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}. \end{aligned}$$

We saw in Chapter 3 that the rank of a symmetric matrix is equal to the number of its nonzero eigenvalues. The diagonal factorization given in (4.12) immediately yields the following generalization of this result in Theorem 4.9.

Theorem 4.9 Let A be an $m \times m$ matrix. If A is diagonalizable, then the rank of A is equal to the number of nonzero eigenvalues of A .

The converse of Theorem 4.9 is not true; that is, a matrix need not be diagonalizable for its rank to equal the number of its nonzero eigenvalues.

Example 4.8 Let A , B , and C be the 2×2 matrices given by

$$A = \begin{bmatrix} 1 & 1 \\ 4 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}.$$

The characteristic equation of A simplifies to $(\lambda - 3)(\lambda + 1) = 0$, so its eigenvalues are $\lambda = 3, -1$. Since the eigenvalues are simple, A is diagonalizable. Eigenvectors corresponding to these two eigenvalues are $\mathbf{x}_1 = (1, 2)'$ and $\mathbf{x}_2 = (1, -2)'$, so the diagonalization of A is given by

$$\begin{bmatrix} 1/2 & 1/4 \\ 1/2 & -1/4 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 4 & 1 \end{bmatrix} \begin{bmatrix} 1 & 1 \\ 2 & -2 \end{bmatrix} = \begin{bmatrix} 3 & 0 \\ 0 & -1 \end{bmatrix}.$$

Clearly, the rank of A is 2, which is the same as the number of nonzero eigenvalues of A . The characteristic equation of B reduces to $\lambda^2 = 0$, so B has the eigenvalue $\lambda = 0$ with multiplicity $r = 2$. Since $\text{rank}(B - \lambda I_2) = \text{rank}(B) = 1 \neq 0 = m - r$, B will not have two linearly independent eigenvectors. The equation $B\mathbf{x} = \lambda\mathbf{x} = \mathbf{0}$ has only one linearly independent solution for \mathbf{x} , namely, vectors of the form $(a, 0)'$. Thus, B is not diagonalizable. Note also that the rank of B is 1, which is greater than the number of its nonzero eigenvalues. Finally, turning to C , we see that it has

the eigenvalue $\lambda = 1$ with multiplicity $r = 2$, because its characteristic equation simplifies to $(1 - \lambda)^2 = 0$. This matrix is not diagonalizable because $\text{rank}(C - \lambda I_2) = \text{rank}(C - I_2) = \text{rank}(B) = 1 \neq 0 = m - r$. Any eigenvector of C is a scalar multiple of the vector $\mathbf{x} = (1, 0)'$. However, notice that even though C is not diagonalizable, it has rank of 2, which is the same as the number of its nonzero eigenvalues.

Theorem 4.10 shows that the connection between the rank and the number of nonzero eigenvalues of a matrix A hinges on the dimension of the eigenspace associated with the eigenvalue 0.

Theorem 4.10 Let A be an $m \times m$ matrix, let k be the dimension of the eigenspace associated with the eigenvalue 0 if 0 is an eigenvalue of A , and let $k = 0$ otherwise. Then

$$\text{rank}(A) = m - k.$$

Proof. From Theorem 2.24, we know that

$$\text{rank}(A) = m - \dim\{N(A)\},$$

where $N(A)$ is the null space of A . However, because the null space of A consists of all vectors \mathbf{x} satisfying $A\mathbf{x} = \mathbf{0}$, we see that $N(A)$ is the same as $S_A(0)$, and so the result follows. \square

We have seen that the number of nonzero eigenvalues of a matrix A equals the rank of A if A is similar to a diagonal matrix; that is, A being diagonalizable is a sufficient condition for this exact relationship between rank and the number of nonzero eigenvalues. The following necessary and sufficient condition for this relationship to exist is an immediate consequence of Theorem 4.10

Corollary 4.10.1 Let A be an $m \times m$ matrix, and let m_0 denote the multiplicity of the eigenvalue 0. Then the rank of A is equal to the number of nonzero eigenvalues of A if and only if

$$\dim\{S_A(0)\} = m_0.$$

Example 4.9 We saw in Example 4.8 that the two matrices

$$B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}, \quad C = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix},$$

are not diagonalizable because each has only one linearly independent eigenvector associated with its single eigenvalue, which has multiplicity 2. This eigenvalue is 0 for B , so

$$\text{rank}(B) = 2 - \dim\{S_B(0)\} = 2 - 1 = 1.$$

On the other hand, because 0 is not an eigenvalue of C , $\dim\{S_C(0)\} = 0$, and so the rank of C equals the number of its nonzero eigenvalues, 2.

4.5 THE JORDAN DECOMPOSITION

Our next factorization of a square matrix A is one that could be described as an attempt to find a matrix similar to A , which, if not diagonal, is as close to being diagonal as is possible. We begin with Definition 4.2.

Definition 4.2 For $h > 1$, the $h \times h$ matrix $J_h(\lambda)$ is said to be a Jordan block matrix if it has the form

$$J_h(\lambda) = \lambda I_h + \sum_{i=1}^{h-1} e_i e'_{i+1} = \begin{bmatrix} \lambda & 1 & 0 & \cdots & 0 \\ 0 & \lambda & 1 & \cdots & 0 \\ 0 & 0 & \lambda & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & \lambda \end{bmatrix},$$

where e_i is the i th column of I_h . If $h = 1$, $J_1(\lambda) = \lambda$.

The matrices B and C from Example 4.8 and Example 4.9 are both 2×2 Jordan block matrices; in particular, $B = J_2(0)$ and $C = J_2(1)$. We saw that neither of these matrices is similar to a diagonal matrix. This is true for Jordan block matrices in general; if $h > 1$, then $J_h(\lambda)$ is not diagonalizable. To see this, note that because $J_h(\lambda)$ is a triangular matrix, its diagonal elements are its eigenvalues, and so it has the one value, λ , repeated h times. However, the solution to $J_h(\lambda)\mathbf{x} = \lambda\mathbf{x}$ has x_1 arbitrary, whereas $x_2 = \cdots = x_h = 0$; that is, $J_h(\lambda)$ has only one linearly independent eigenvector, which is of the form $\mathbf{x} = (x_1, 0, \dots, 0)'$.

We now state the Jordan decomposition theorem. For a proof of this result, see Horn and Johnson (2013).

Theorem 4.11 Let A be an $m \times m$ matrix. Then a nonsingular matrix B exists, such that

$$\begin{aligned} B^{-1}AB &= J = \text{diag}(J_{h_1}(\lambda_1), \dots, J_{h_r}(\lambda_r)) \\ &= \begin{bmatrix} J_{h_1}(\lambda_1) & (0) & \cdots & (0) \\ (0) & J_{h_2}(\lambda_2) & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & J_{h_r}(\lambda_r) \end{bmatrix}, \end{aligned}$$

where $h_1 + \cdots + h_r = m$ and $\lambda_1, \dots, \lambda_r$ are the not necessarily distinct eigenvalues of A .

The matrix J in Theorem 4.11 will be diagonal if $h_i = 1$ for all i . Since the $h_i \times h_i$ matrix $J_{h_i}(\lambda_i)$ has only one linearly independent eigenvector, it follows that the Jordan canonical form $J = \text{diag}(J_{h_1}(\lambda_1), \dots, J_{h_r}(\lambda_r))$ has r linearly independent eigenvectors. Thus, if $h_i > 1$ for at least one i , then J will not be diagonal; in fact, J will not be diagonalizable. The vector \mathbf{x}_i is an eigenvector of J corresponding to the eigenvalue λ_i if and only if the vector $\mathbf{y}_i = B\mathbf{x}_i$ is an eigenvector of A corresponding to λ_i ; for instance, if \mathbf{x}_i satisfies $J\mathbf{x}_i = \lambda_i\mathbf{x}_i$, then

$$A\mathbf{y}_i = (BJB^{-1})B\mathbf{x}_i = BJ\mathbf{x}_i = \lambda_i B\mathbf{x}_i = \lambda_i \mathbf{y}_i.$$

Thus, r also gives the number of linearly independent eigenvectors of A , and A is diagonalizable only if J is diagonal.

Example 4.10 Suppose that A is a 4×4 matrix with the eigenvalue λ having multiplicity 4. Then A will be similar to one of the following five Jordan canonical forms:

$$\begin{aligned} \text{diag}(J_1(\lambda), J_1(\lambda), J_1(\lambda), J_1(\lambda)) &= \begin{bmatrix} \lambda & 0 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix}, \\ \text{diag}(J_2(\lambda), J_1(\lambda), J_1(\lambda)) &= \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix}, \\ \text{diag}(J_3(\lambda), J_1(\lambda)) &= \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 1 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix}, \\ \text{diag}(J_2(\lambda), J_2(\lambda)) &= \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 1 \\ 0 & 0 & 0 & \lambda \end{bmatrix}, \\ J_4(\lambda) &= \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 1 & 0 \\ 0 & 0 & \lambda & 1 \\ 0 & 0 & 0 & \lambda \end{bmatrix}. \end{aligned}$$

The first form given is diagonal, so this corresponds to the case in which A has four linearly independent eigenvectors associated with the eigenvalue λ . The second and last forms correspond to A having three and one linearly independent eigenvectors, respectively. If A has two linearly independent eigenvectors, then it will be similar to either the third or the fourth matrix given.

4.6 THE SCHUR DECOMPOSITION

Our next result can be viewed as another generalization of the spectral decomposition theorem to any square matrix A . The diagonalization theorem and the Jordan decomposition were generalizations of the spectral decomposition in which our goal was to obtain a diagonal or “nearly” diagonal matrix. Now, instead we focus on the orthogonal matrix employed in the spectral decomposition theorem. Specifically, if we restrict attention only to orthogonal matrices, X , what is the simplest structure that we can get for $X'AX$? It turns out that for the general case of any real square matrix A , we can find an X such that X^*AX is a triangular matrix, where we have broadened the choice of X to include all unitary matrices. Recall that a real unitary matrix is an orthogonal matrix, and in general, X is unitary if $X^*X = I$, where X^* is the transpose of the complex conjugate of X . This decomposition, sometimes referred to as the Schur decomposition, is given in Theorem 4.12.

Theorem 4.12 Let A be an $m \times m$ matrix. Then an $m \times m$ unitary matrix X exists, such that

$$X^*AX = T,$$

where T is an upper triangular matrix with the eigenvalues of A as its diagonal elements.

Proof. Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of A , and let \mathbf{y}_1 be an eigenvector of A corresponding to λ_1 and normalized so that $\mathbf{y}_1^*\mathbf{y}_1 = 1$. Let Y be any $m \times m$ unitary matrix having \mathbf{y}_1 as its first column. Writing Y in partitioned form as $Y = [\mathbf{y}_1 \ Y_2]$, we see that, because $A\mathbf{y}_1 = \lambda_1\mathbf{y}_1$ and $Y_2^*\mathbf{y}_1 = \mathbf{0}$,

$$\begin{aligned} Y^*AY &= \begin{bmatrix} \mathbf{y}_1^*A\mathbf{y}_1 & \mathbf{y}_1^*AY_2 \\ Y_2^*A\mathbf{y}_1 & Y_2^*AY_2 \end{bmatrix} = \begin{bmatrix} \lambda_1\mathbf{y}_1^*\mathbf{y}_1 & \mathbf{y}_1^*AY_2 \\ \lambda_1Y_2^*\mathbf{y}_1 & Y_2^*AY_2 \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & \mathbf{y}_1^*AY_2 \\ \mathbf{0} & B \end{bmatrix}, \end{aligned}$$

where the $(m-1) \times (m-1)$ matrix $B = Y_2^*AY_2$. Using the identity above and the cofactor expansion formula for a determinant, it follows that the characteristic equation of Y^*AY is

$$(\lambda_1 - \lambda)|B - \lambda I_{m-1}| = 0,$$

and, because by Theorem 3.2(d) the eigenvalues of Y^*AY are the same as those of A , the eigenvalues of B must be $\lambda_2, \dots, \lambda_m$. Now if $m = 2$, then the scalar B must equal λ_2 and Y^*AY is upper triangular, so the proof is complete. For $m > 2$, we proceed by induction; that is, we show that if our result holds for $(m-1) \times (m-1)$ matrices, then it must also hold for $m \times m$ matrices. Since B is $(m-1) \times (m-1)$, we may assume that a unitary matrix W exists, such that $W^*BW = T_2$, where T_2 is

an upper triangular matrix with diagonal elements $\lambda_2, \dots, \lambda_m$. Define the $m \times m$ matrix U by

$$U = \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W \end{bmatrix},$$

and note that U is unitary because W is. If we let $X = YU$, then X is also unitary and

$$\begin{aligned} X^*AX &= U^*Y^*AYU = \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W^* \end{bmatrix} \begin{bmatrix} \lambda_1 & \mathbf{y}_1^*AY_2 \\ \mathbf{0} & B \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & \mathbf{y}_1^*AY_2W \\ \mathbf{0} & W^*BW \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & \mathbf{y}_1^*AY_2W \\ \mathbf{0} & T_2 \end{bmatrix}, \end{aligned}$$

where this final matrix is upper triangular with $\lambda_1, \dots, \lambda_m$ as its diagonal elements. Thus, the proof is complete. \square

If all of the eigenvalues of A are real, then corresponding real eigenvectors exist. In this case, a real matrix X satisfying the conditions of Theorem 4.12 can be found. Consequently, we have the following result.

Corollary 4.12.1 If the $m \times m$ matrix A has real eigenvalues, then an $m \times m$ orthogonal matrix X exists, such that $X'AX = T$, where T is an upper triangular matrix.

Example 4.11 Consider the 3×3 matrix given by

$$A = \begin{bmatrix} 5 & -3 & 3 \\ 4 & -2 & 3 \\ 4 & -4 & 5 \end{bmatrix}.$$

In Example 3.1, the eigenvalues of A were shown to be $\lambda_1 = 1$, $\lambda_2 = 2$, and $\lambda_3 = 5$, with eigenvectors, $\mathbf{x}_1 = (0, 1, 1)'$, $\mathbf{x}_2 = (1, 1, 0)'$, and $\mathbf{x}_3 = (1, 1, 1)'$, respectively. We will find an orthogonal matrix X and an upper triangular matrix T so that $A = XTX'$. First, we construct an orthogonal matrix Y having a normalized version of \mathbf{x}_1 as its first column; for instance, by inspection, we set

$$Y = \begin{bmatrix} 0 & 0 & 1 \\ 1/\sqrt{2} & 1/\sqrt{2} & 0 \\ 1/\sqrt{2} & -1/\sqrt{2} & 0 \end{bmatrix}.$$

Thus, our first stage yields

$$Y'AY = \begin{bmatrix} 1 & -7 & 4\sqrt{2} \\ 0 & 2 & 0 \\ 0 & -3\sqrt{2} & 5 \end{bmatrix}.$$

The 2×2 matrix

$$B = \begin{bmatrix} 2 & 0 \\ -3\sqrt{2} & 5 \end{bmatrix}$$

has a normalized eigenvector $(1/\sqrt{3}, \sqrt{2}/\sqrt{3})'$, and so we can construct an orthogonal matrix

$$W = \begin{bmatrix} 1/\sqrt{3} & -\sqrt{2}/\sqrt{3} \\ \sqrt{2}/\sqrt{3} & 1/\sqrt{3} \end{bmatrix}$$

for which

$$W'BW = \begin{bmatrix} 2 & 3\sqrt{2} \\ 0 & 5 \end{bmatrix}.$$

Putting it all together, we have

$$X = Y \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W \end{bmatrix} = \frac{1}{\sqrt{6}} \begin{bmatrix} 0 & 2 & \sqrt{2} \\ \sqrt{3} & 1 & -\sqrt{2} \\ \sqrt{3} & -1 & \sqrt{2} \end{bmatrix}$$

and

$$T = X'AX = \begin{bmatrix} 1 & 1/\sqrt{3} & 22/\sqrt{6} \\ 0 & 2 & 3\sqrt{2} \\ 0 & 0 & 5 \end{bmatrix}.$$

The matrices X and T in the Schur decomposition are not unique; that is, if $A = XTX^*$ is a Schur decomposition of A , then $A = X_0T_0X_0^*$ is also, where $X_0 = XP$ and P is any unitary matrix for which $P^*TP = T_0$ is upper triangular. The triangular matrices T and T_0 must have the same diagonal elements, possibly ordered differently. Otherwise, however, the two matrices T and T_0 may be quite different. For example, it can be easily verified that the matrices

$$X_0 = \begin{bmatrix} 1/\sqrt{3} & 2/\sqrt{6} & 0 \\ 1/\sqrt{3} & -1/\sqrt{6} & -1/\sqrt{2} \\ 1/\sqrt{3} & -1/\sqrt{6} & 1/\sqrt{2} \end{bmatrix}, \quad T_0 = \begin{bmatrix} 5 & 8/\sqrt{2} & 20/\sqrt{6} \\ 0 & 1 & -1/\sqrt{3} \\ 0 & 0 & 2 \end{bmatrix}$$

give another Schur decomposition of the matrix A .

In Chapter 3, by using the characteristic equation of the $m \times m$ matrix A , we were able to prove that the determinant of A equals the product of its eigenvalues, whereas the trace of A equals the sum of its eigenvalues. These results are also very easily proven using the Schur decomposition of A . If the eigenvalues of A are $\lambda_1, \dots, \lambda_m$ and $A = XTX^*$ is a Schur decomposition of A , then it follows that

$$|A| = |XTX^*| = |X^*X||T| = |T| = \prod_{i=1}^m \lambda_i,$$

because $|X^*X| = 1$ follows from the fact that X is a unitary matrix, and the determinant of a triangular matrix is the product of its diagonal elements. Also, using properties of the trace of a matrix, we have

$$\text{tr}(A) = \text{tr}(XTX^*) = \text{tr}(X^*XT) = \text{tr}(T) = \sum_{i=1}^m \lambda_i.$$

The Schur decomposition also provides a method of easily establishing the fact that the number of nonzero eigenvalues of a matrix serves as a lower bound for the rank of that matrix. This is the subject of our next theorem.

Theorem 4.13 Suppose the $m \times m$ matrix A has r nonzero eigenvalues. Then $\text{rank}(A) \geq r$.

Proof. Let X be a unitary matrix and T be an upper triangular matrix, such that $A = XTX^*$. Since the eigenvalues of A are the diagonal elements of T , T must have exactly r nonzero diagonal elements. The $r \times r$ submatrix of T , formed by deleting the columns and rows occupied by the zero diagonal elements of T , will be upper triangular with nonzero diagonal elements. This submatrix will be nonsingular because the determinant of a triangular matrix is the product of its diagonal elements, so we must have $\text{rank}(T) \geq r$. However, because X is unitary, it must be nonsingular, so

$$\text{rank}(A) = \text{rank}(XTX^*) = \text{rank}(T) \geq r,$$

and the proof is complete. □

4.7 THE SIMULTANEOUS DIAGONALIZATION OF TWO SYMMETRIC MATRICES

We have already discussed in Section 3.8 one manner in which two symmetric matrices can be simultaneously diagonalized. We restate this result in Theorem 4.14.

Theorem 4.14 Let A and B be $m \times m$ symmetric matrices with B being positive definite. Let $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$, where $\lambda_1, \dots, \lambda_m$ are the eigenvalues of $B^{-1}A$. Then a nonsingular matrix C exists, such that

$$CAC' = \Lambda, \quad CBC' = I_m.$$

Example 4.12 One application of the simultaneous diagonalization described in Theorem 4.14 is in a multivariate analysis commonly referred to as canonical variate analysis (see Krzanowski, 2000 or Mardia, et al. 1979). This analysis involves data from the multivariate one-way classification model discussed in Example 3.16, so that we have independent random samples from k different groups or treatments, with the

i th sample of $m \times 1$ vectors given by $\mathbf{y}_{i1}, \dots, \mathbf{y}_{in_i}$. The model is

$$\mathbf{y}_{ij} = \boldsymbol{\mu}_i + \boldsymbol{\epsilon}_{ij},$$

where $\boldsymbol{\mu}_i$ is an $m \times 1$ vector of constants and $\boldsymbol{\epsilon}_{ij} \sim N_m(\mathbf{0}, \Omega)$. In Example 3.16, we saw how the matrices

$$B = \sum_{i=1}^k n_i (\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)',$$

where

$$\bar{\mathbf{y}}_i = \sum_{j=1}^{n_i} \mathbf{y}_{ij}/n_i, \quad \bar{\mathbf{y}} = \sum_{i=1}^k n_i \bar{\mathbf{y}}_i/n, \quad n = \sum_{i=1}^k n_i,$$

could be used to test the hypothesis, $H_0 : \boldsymbol{\mu}_1 = \dots = \boldsymbol{\mu}_k$. Canonical variate analysis is an analysis of the differences in the mean vectors, performed when this hypothesis is rejected. This analysis is particularly useful when the differences between the vectors $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ are confined, or nearly confined, to some lower dimensional subspace of R^m . Note that if these vectors span an r -dimensional subspace of R^m , then the population version of B ,

$$\Phi = \sum_{i=1}^k n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})',$$

where $\boldsymbol{\mu} = \sum n_i \boldsymbol{\mu}_i/n$, will have rank r ; in fact, the eigenvectors of Φ corresponding to its positive eigenvalues will span this r -dimensional space. Thus, a plot of the projections of $\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_k$ onto this subspace will yield a reduced-dimension diagram of the population means. Unfortunately, if $\Omega \neq I_m$, it will be difficult to interpret the differences in these mean vectors because Euclidean distance would not be appropriate. This difficulty can be resolved by analyzing the transformed data $\Omega^{-1/2} \mathbf{y}_{ij}$, where $\Omega^{-1/2} \Omega^{-1/2} = \Omega^{-1}$, because $\Omega^{-1/2} \mathbf{y}_{ij} \sim N_m(\Omega^{-1/2} \boldsymbol{\mu}_i, I_m)$. Thus, we would plot the projections of $\Omega^{-1/2} \boldsymbol{\mu}_1, \dots, \Omega^{-1/2} \boldsymbol{\mu}_k$ onto the subspace spanned by the eigenvectors of $\Omega^{-1/2} \Phi \Omega^{-1/2}$ corresponding to its r positive eigenvalues; that is, if the spectral decomposition of $\Omega^{-1/2} \Phi \Omega^{-1/2}$ is given by $P_1 \Lambda_1 P_1'$, where P_1 is an $m \times r$ matrix satisfying $P_1' P_1 = I_r$ and Λ_1 is an $r \times r$ diagonal matrix, then we could simply plot the vectors $P_1' \Omega^{-1/2} \boldsymbol{\mu}_1, \dots, P_1' \Omega^{-1/2} \boldsymbol{\mu}_k$ in R^r . The r components of the vector $\mathbf{v}_i = P_1' \Omega^{-1/2} \boldsymbol{\mu}_i$ in this r -dimensional space are called the canonical variates means for the i th population. Note that in obtaining these canonical variates, we have essentially used the simultaneous diagonalization of Φ and Ω , because if $C' = (C'_1, C'_2)$ satisfies

$$\begin{bmatrix} C_1 \\ C_2 \end{bmatrix} \Phi \begin{bmatrix} C'_1 & C'_2 \end{bmatrix} = \begin{bmatrix} \Lambda_1 & (0) \\ (0) & (0) \end{bmatrix},$$

$$\begin{bmatrix} C_1 \\ C_2 \end{bmatrix} \Omega [C'_1 \quad C'_2] = \begin{bmatrix} I_r & (0) \\ (0) & I_{m-r} \end{bmatrix},$$

then we can take $C_1 = P'_1 \Omega^{-1/2}$. When μ_1, \dots, μ_k are unknown, the canonical variate means can be estimated by the sample canonical variate means, which are computed using the sample means $\bar{y}_1, \dots, \bar{y}_k$ and the corresponding simultaneous diagonalization of B and W .

Theorem 4.14 is a special case of the more general result given next.

Theorem 4.15 Let A and B be $m \times m$ symmetric matrices, and suppose that there is a linear combination of A and B that is positive definite. Then a nonsingular matrix C exists, such that both CAC' and CBC' are diagonal.

Proof. Let $D = \alpha A + \beta B$ be a linear combination of A and B for which D is positive definite. If both α and β are zero, D would not be positive definite, so we may assume without loss of generality that $\alpha \neq 0$. In this case, A may be written as $A = \alpha^{-1}(D - \beta B)$. Since D is positive definite, a nonsingular matrix T exists, such that $D = TT'$, or equivalently, $T^{-1}DT^{-1'} = I_m$. Further, $T^{-1}BT^{-1'}$ is symmetric, so an orthogonal matrix P exists, for which $P'T^{-1}BT^{-1'}P = \Delta$ is diagonal. Thus, if we define $C = P'T^{-1}$, we have $CDC' = P'P = I_m$ and $CBC' = \Delta$; that is, B is diagonalized by C and so also is A because $CAC' = \alpha^{-1}(CDC' - \beta CBC') = \alpha(I_m - \beta\Delta)$. \square

We can get another set of sufficient conditions for A and B to be simultaneously diagonalizable by strengthening the condition on A given in Theorem 4.14 while weakening the condition on B .

Theorem 4.16 Let A and B be $m \times m$ nonnegative definite matrices. Then a nonsingular matrix C exists, such that both CAC' and CBC' are diagonal.

Proof. Let $r_1 = \text{rank}(A)$, $r_2 = \text{rank}(B)$, and assume without loss of generality that $r_1 \leq r_2$. Let $A = P_1 \Lambda_1 P'_1$ be the spectral decomposition of A , so that the $m \times r_1$ matrix P_1 satisfies $P'_1 P_1 = I_{r_1}$ and Λ_1 is an $r_1 \times r_1$ diagonal matrix with positive diagonal elements. Define $A_1 = C_1 A C'_1$ and $B_1 = C_1 B C'_1$, where $C'_1 = [P_1 \Lambda_1^{-1/2} \quad P_2]$ and P_2 is any $m \times (m - r_1)$ matrix for which $[P_1 \quad P_2]$ is orthogonal, and note that

$$A_1 = \begin{bmatrix} I_{r_1} & (0) \\ (0) & (0) \end{bmatrix},$$

so that C_1 diagonalizes A . If any of the last $m - r_1$ diagonal elements of B_1 is zero, then all elements in that row and column are zero because B_1 is nonnegative definite (Problem 3.51). Note that if the $(r_1 + i, r_1 + i)$ th element of B_1 is $b \neq 0$ and the $(r_1 + i, j)$ th and $(j, r_1 + i)$ th elements of B_1 are a , and we define

$$T = I_m - \frac{a}{b} e_j e'_{r_1+i},$$

then TB_1T' yields a matrix identical to B_1 except that a multiple of row $r_1 + i$ has been added to row j after which a multiple of column $r_1 + i$ has been added to column j so that the $(r_1 + i, j)$ th and $(j, r_1 + i)$ th elements are now each zero. We can repeatedly use this process to get a nonsingular matrix C_2 , which is a product of matrices of the form given for T so that

$$C_2B_1C_2' = \begin{bmatrix} B_* & (0) \\ (0) & D_1 \end{bmatrix},$$

where B_* is an $r_1 \times r_1$ matrix of rank r_3 , D_1 is an $(m - r_1) \times (m - r_1)$ diagonal matrix with r_4 nonzero diagonal elements each of which is positive, and $r_3 + r_4 = r_2$. Since the matrix T , and hence also the matrix C_2 , when partitioned has the form

$$\begin{bmatrix} I_{r_1} & E \\ (0) & F \end{bmatrix}$$

for some $r_1 \times (m - r_1)$ and $(m - r_1) \times (m - r_1)$ matrices E and F , it follows that $C_2A_1C_2' = A_1$. Finally, define C_3 as

$$C_3 = \begin{bmatrix} Q' & (0) \\ (0) & I_{m-r_1} \end{bmatrix},$$

where Q is an $r_1 \times r_1$ orthogonal matrix satisfying $B_* = QD_2Q'$ and D_2 is an $r_1 \times r_1$ diagonal matrix with r_3 nonzero diagonal elements each of which is positive. Then with $C = C_3C_2C_1$, we have $CAC' = \text{diag}(I_{r_1}, (0))$ and $CBC' = \text{diag}(D_2, D_1)$. \square

We are now in a position to establish a useful determinantal inequality.

Theorem 4.17 Suppose A and B are $m \times m$ nonnegative definite matrices. Then

$$|A + B| \geq |A| + |B|,$$

with equality if and only if $A = (0)$ or $B = (0)$ or $A + B$ is singular.

Proof. Since $A + B$ is also nonnegative definite, the inequality clearly holds when $|A| = |B| = 0$, with equality if and only if $A + B$ is singular. For the remainder of the proof we assume, without loss of generality, that B is positive definite. Using Theorem 4.14, we find that we can establish the result by showing that

$$\prod_{i=1}^m (\lambda_i + 1) = |\Lambda + I_m| \geq |\Lambda| + |I_m| = \prod_{i=1}^m \lambda_i + 1, \quad (4.13)$$

with equality if and only if $\Lambda = (0)$. We prove this by induction. For $m = 2$,

$$(\lambda_1 + 1)(\lambda_2 + 1) = \lambda_1\lambda_2 + \lambda_1 + \lambda_2 + 1 \geq \lambda_1\lambda_2 + 1$$

since λ_1 and λ_2 are nonnegative, and we have equality if and only if $\lambda_1 = \lambda_2 = 0$. Now if (4.13) holds for $m - 1$, with equality if and only if $\lambda_1 = \cdots = \lambda_{m-1} = 0$, then

$$\begin{aligned} \prod_{i=1}^m (\lambda_i + 1) &= \left\{ \prod_{i=1}^{m-1} (\lambda_i + 1) \right\} (\lambda_m + 1) \geq \left\{ \prod_{i=1}^{m-1} \lambda_i + 1 \right\} (\lambda_m + 1) \\ &= \prod_{i=1}^m \lambda_i + \prod_{i=1}^{m-1} \lambda_i + \lambda_m + 1 \geq \prod_{i=1}^m \lambda_i + 1, \end{aligned}$$

as is required. The first inequality is an equality if and only if $\lambda_1 = \cdots = \lambda_{m-1} = 0$, while the second is an equality if and only if $\prod_{i=1}^{m-1} \lambda_i = 0$ and $\lambda_m = 0$, and so the proof is complete. \square

The matrix C that diagonalizes A and B in Theorems 4.14, 4.15, and 4.16 is non-singular but not necessarily orthogonal. Further, the diagonal elements of the two diagonal matrices are not the eigenvalues of A nor B . This sort of diagonalization, one which will be useful in our study of quadratic forms in normal random vectors in Chapter 11, is what we consider next; we would like to know whether there exists an orthogonal matrix that diagonalizes both A and B . Theorem 4.18 gives a necessary and sufficient condition for such an orthogonal matrix to exist.

Theorem 4.18 Suppose that A and B are $m \times m$ symmetric matrices. Then an orthogonal matrix P exists, such that $P'AP$ and $P'BP$ are both diagonal if and only if A and B commute; that is, if and only if $AB = BA$.

Proof. First suppose that such an orthogonal matrix does exist; that is, there is an orthogonal matrix P such that $P'AP = \Lambda_1$ and $P'BP = \Lambda_2$, where Λ_1 and Λ_2 are diagonal matrices. Then because Λ_1 and Λ_2 are diagonal matrices, clearly $\Lambda_1\Lambda_2 = \Lambda_2\Lambda_1$, so we have

$$\begin{aligned} AB &= P\Lambda_1P'P\Lambda_2P' = P\Lambda_1\Lambda_2P' = P\Lambda_2\Lambda_1P' \\ &= P\Lambda_2P'P\Lambda_1P' = BA, \end{aligned}$$

and hence, A and B do commute. Conversely, now assuming that $AB = BA$, we need to show that such an orthogonal matrix P does exist. Let μ_1, \dots, μ_h be the distinct values of the eigenvalues of A having multiplicities r_1, \dots, r_h , respectively. Since A is symmetric, an orthogonal matrix Q exists, satisfying

$$Q'AQ = \Lambda_1 = \text{diag}(\mu_1 I_{r_1}, \dots, \mu_h I_{r_h}).$$

Performing this same transformation on B and partitioning the resulting matrix in the same way that $Q'AQ$ has been partitioned, we get

$$C = Q'BQ = \begin{bmatrix} C_{11} & C_{12} & \cdots & C_{1h} \\ C_{21} & C_{22} & \cdots & C_{2h} \\ \vdots & \vdots & & \vdots \\ C_{h1} & C_{h2} & \cdots & C_{hh} \end{bmatrix},$$

where C_{ij} is $r_i \times r_j$. Note that because $AB = BA$, we must have

$$\begin{aligned} \Lambda_1 C &= Q'AQQ'BQ = Q'ABQ = Q'BAQ \\ &= Q'BQQ'AQ = C\Lambda_1. \end{aligned}$$

Equating the (i, j) th submatrix of $\Lambda_1 C$ to the (i, j) th submatrix of $C\Lambda_1$ yields the identity $\mu_i C_{ij} = \mu_j C_{ij}$. Since $\mu_i \neq \mu_j$ if $i \neq j$, we must have $C_{ij} = (0)$ if $i \neq j$; that is, the matrix $C = \text{diag}(C_{11}, \dots, C_{hh})$ is block diagonal. Now because C is symmetric, so also is C_{ii} for each i , and thus, we can find an $r_i \times r_i$ orthogonal matrix X_i satisfying

$$X_i' C_{ii} X_i = \Delta_i,$$

where Δ_i is diagonal. Let $P = QX$, where X is the block diagonal matrix $X = \text{diag}(X_1, \dots, X_h)$, and note that

$$\begin{aligned} P'P &= X'Q'QX = X'X \\ &= \text{diag}(X_1'X_1, \dots, X_h'X_h) \\ &= \text{diag}(I_{r_1}, \dots, I_{r_h}) = I_m, \end{aligned}$$

so that P is orthogonal. Finally, the matrix $\Delta = \text{diag}(\Delta_1, \dots, \Delta_h)$ is diagonal and

$$\begin{aligned} P'AP &= X'Q'AQX = X'\Lambda_1X \\ &= \text{diag}(X_1', \dots, X_h') \text{diag}(\mu_1 I_{r_1}, \dots, \mu_h I_{r_h}) \text{diag}(X_1, \dots, X_h) \\ &= \text{diag}(\mu_1 X_1'X_1, \dots, \mu_h X_h'X_h) \\ &= \text{diag}(\mu_1 I_{r_1}, \dots, \mu_h I_{r_h}) = \Lambda_1 \end{aligned}$$

and

$$\begin{aligned} P'BP &= X'Q'BQX = X'CX \\ &= \text{diag}(X_1', \dots, X_h') \text{diag}(C_{11}, \dots, C_{hh}) \text{diag}(X_1, \dots, X_h) \\ &= \text{diag}(X_1' C_{11} X_1, \dots, X_h' C_{hh} X_h) \\ &= \text{diag}(\Delta_1, \dots, \Delta_h) = \Delta, \end{aligned}$$

and so the proof is complete. \square

The columns of the matrix P are eigenvectors of A as well as B ; that is, A and B commute if and only if the two matrices have common eigenvectors. Also, note that because A and B are symmetric, $(AB)' = B'A' = BA$, and so $AB = BA$ if and only if AB is symmetric. Theorem 4.18 easily generalizes to a collection of symmetric matrices.

Theorem 4.19 Let A_1, \dots, A_k be $m \times m$ symmetric matrices. Then an orthogonal matrix P exists, such that $P'A_iP = \Lambda_i$ is diagonal for each i if and only if $A_iA_j = A_jA_i$ for all pairs (i, j) .

The two previous theorems involving symmetric matrices are special cases of more general results regarding diagonalizable matrices. For instance, Theorem 4.19 is a special case of the following result. The proof, which is similar to that given for Theorem 4.18, is left as an exercise.

Theorem 4.20 Suppose that each of the $m \times m$ matrices A_1, \dots, A_k is diagonalizable. Then a nonsingular matrix X exists, such that $X^{-1}A_iX = \Lambda_i$ is diagonal for each i if and only if $A_iA_j = A_jA_i$ for all pairs (i, j) .

4.8 MATRIX NORMS

In Chapter 2, we saw that vector norms can be used to measure the size of a vector. Similarly, we may be interested in measuring the size of an $m \times m$ matrix A or measuring the closeness of A to another $m \times m$ matrix B . Matrix norms will provide the means to do this. In a later chapter, we will need to apply some of our results on matrix norms to matrices that are possibly complex matrices. Consequently, throughout this section, we will not be restricting attention only to real matrices.

Definition 4.3 A function $\|A\|$ defined on all $m \times m$ matrices A , real or complex, is a matrix norm if the following conditions hold for all $m \times m$ matrices A and B :

- (a) $\|A\| \geq 0$.
- (b) $\|A\| = 0$ if and only if $A = (0)$.
- (c) $\|cA\| = |c| \|A\|$ for any complex scalar c .
- (d) $\|A + B\| \leq \|A\| + \|B\|$.
- (e) $\|AB\| \leq \|A\| \|B\|$.

Any vector norm defined on $m^2 \times 1$ vectors, when applied to the $m^2 \times 1$ vector formed by stacking the columns of A , one on top of the other, will satisfy conditions (a)–(d) because these are the conditions of a vector norm. However, condition (e), which relates the sizes of A and B to that of AB , will not necessarily hold for vector norms; that is, not all vector norms can be used as matrix norms.

Although a given vector norm might not be a matrix norm, it can always be used to find a matrix norm.

Theorem 4.21 Suppose $\|\mathbf{x}\|$ is a vector norm defined on $m \times 1$ vectors. Then

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

is a matrix norm defined on $m \times m$ matrices.

Proof. We simply need to verify the five conditions of Definition 4.3. Condition (a) follows immediately since the nonnegativity of $\|A\|$ is guaranteed by the nonnegativity of $\|A\mathbf{x}\|$. Clearly, $A = (0)$ implies $\|A\| = 0$, while $\|A\| = 0$ implies $\|A\mathbf{e}_i\| = 0$, that is, $A\mathbf{e}_i = \mathbf{0}$, for $i = 1, \dots, m$, so (b) holds. Condition (c) holds since for any scalar c

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

Next, for any unit vector \mathbf{x} , we have

$$\|(A + B)\mathbf{x}\| = \|A\mathbf{x} + B\mathbf{x}\| \leq \|A\mathbf{x}\| + \|B\mathbf{x}\| \leq \|A\| + \|B\|,$$

and so

$$\|A + B\| = \max_{\|\mathbf{x}\|=1} \|(A + B)\mathbf{x}\| \leq \|A\| + \|B\|.$$

To establish (e), we first note that $\|A\mathbf{y}\| \leq \|A\|\|\mathbf{y}\|$ holds, clearly when $\mathbf{y} = \mathbf{0}$, and when $\mathbf{y} \neq \mathbf{0}$ since

$$\|A\| = \max_{\|\mathbf{x}\|=1} \|A\mathbf{x}\| \geq \|A\| \frac{\|\mathbf{y}\|}{\|\mathbf{y}\|} = \|A\mathbf{y}\| / \|\mathbf{y}\|$$

yields the required inequality. Thus, for any unit vector \mathbf{x} ,

$$\|AB\mathbf{x}\| = \|A(B\mathbf{x})\| \leq \|A\|\|B\mathbf{x}\| \leq \|A\|\|B\|,$$

so

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

□

The matrix norm, $\|A\|$, given in Theorem 4.21 is sometimes referred to as the matrix norm induced by the vector norm $\|\mathbf{x}\|$.

We now give examples of some commonly encountered matrix norms. We will leave it to the reader to verify that these functions, in fact, satisfy the conditions

of Definition 4.3. The Euclidean matrix norm is simply the Euclidean vector norm computed on the stacked columns of A , and so it is given by

$$\|A\|_E = \left(\sum_{i=1}^m \sum_{j=1}^m |a_{ij}|^2 \right)^{1/2} = \{\text{tr}(A^*A)\}^{1/2}.$$

The maximum column sum matrix norm is given by

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

whereas the maximum row sum matrix norm is given by

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

The spectral norm uses the eigenvalues of A^*A ; in particular, if μ_1, \dots, μ_m are the eigenvalues of A^*A , then the spectral norm is given by

$$\|A\|_2 = \max_{1 \leq i \leq m} \sqrt{\mu_i};$$

that is, the spectral norm of A is the largest singular value of A .

Example 4.13 In this example, we will show that the spectral matrix norm $\|A\|_2 = \sqrt{\mu_1}$, where $\mu_1 \geq \dots \geq \mu_m \geq 0$ are the eigenvalues of A^*A , is induced by the Euclidean vector norm $\|x\|_2 = (x^*x)^{1/2}$. Let $A = PDQ^*$ be the singular value decomposition of A so that P and Q are $m \times m$ unitary matrices and $D = \text{diag}(\sqrt{\mu_1}, \dots, \sqrt{\mu_m})$. Put $y = Q^*x$ so that $x = Qy$ and $\|x\|_2 = \|Qy\|_2 = (y^*Q^*Qy)^{1/2} = (y^*y)^{1/2} = \|y\|_2$. Then

$$\begin{aligned} \max_{\|x\|_2=1} \|Ax\|_2 &= \max_{\|x\|_2=1} \|PDQ^*x\|_2 = \max_{\|x\|_2=1} \{(PDQ^*x)^*(PDQ^*x)\}^{1/2} \\ &= \max_{\|x\|_2=1} (x^*QD^2Q^*x)^{1/2} = \max_{\|Qy\|_2=1} (y^*D^2y)^{1/2} \\ &= \max_{\|y\|_2=1} (y^*D^2y)^{1/2} = \max_{\|y\|_2=1} \left(\sum_{i=1}^m \mu_i |y_i|^2 \right)^{1/2} \\ &\leq \sqrt{\mu_1} \max_{\|y\|_2=1} \left(\sum_{i=1}^m |y_i|^2 \right)^{1/2} = \sqrt{\mu_1}. \end{aligned}$$

The inequality can be replaced by equality since if q_1 is the first column of Q , $\|Aq_1\|_2 = \sqrt{\mu_1}$ and $\|q_1\|_2 = 1$.

We will find the following theorem useful. The proof, which simply involves the verification of the conditions of Definition 4.3, is left to the reader as an exercise.

Theorem 4.22 Let $\|A\|$ be any matrix norm defined on $m \times m$ matrices. If C is an $m \times m$ nonsingular matrix, then the function defined by

$$\|A\|_C = \|C^{-1}AC\|$$

is also a matrix norm.

The eigenvalues of a matrix A play an important role in the study of matrix norms of A . Particularly important is the maximum modulus of this set of eigenvalues.

Definition 4.4 Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of the $m \times m$ matrix A . The spectral radius of A , denoted $\rho(A)$, is defined to be

$$\rho(A) = \max_{1 \leq i \leq m} |\lambda_i|.$$

Although $\rho(A)$ does give us some information about the size of A , it is not a matrix norm itself. To see this, consider the case in which $m = 2$ and

$$A = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

Both of the eigenvalues of A are 0, so $\rho(A) = 0$ even though A is not the null matrix; that is, $\rho(A)$ violates condition (b) of Definition 4.3. Theorem 4.23 shows that $\rho(A)$ actually serves as a lower bound for any matrix norm of A .

Theorem 4.23 For any $m \times m$ matrix A and any matrix norm $\|A\|$, $\rho(A) \leq \|A\|$.

Proof. Suppose that λ is an eigenvalue of A for which $|\lambda| = \rho(A)$, and let \mathbf{x} be a corresponding eigenvector, so that $A\mathbf{x} = \lambda\mathbf{x}$. Then $\mathbf{x}\mathbf{1}'_m$ is an $m \times m$ matrix satisfying $A\mathbf{x}\mathbf{1}'_m = \lambda\mathbf{x}\mathbf{1}'_m$, and so using properties (c) and (e) of matrix norms, we find that

$$\rho(A)\|\mathbf{x}\mathbf{1}'_m\| = |\lambda| \|\mathbf{x}\mathbf{1}'_m\| = \|\lambda\mathbf{x}\mathbf{1}'_m\| = \|A\mathbf{x}\mathbf{1}'_m\| \leq \|A\| \|\mathbf{x}\mathbf{1}'_m\|.$$

The result now follows by dividing the equation above by $\|\mathbf{x}\mathbf{1}'_m\|$. □

Although the spectral radius of A is at least as small as every norm of A , our next result shows that we can always find a matrix norm so that $\|A\|$ is arbitrarily close to $\rho(A)$.

Theorem 4.24 For any $m \times m$ matrix A and any scalar $\epsilon > 0$, there exists a matrix norm, $\|A\|_{A,\epsilon}$, such that

$$\|A\|_{A,\epsilon} - \rho(A) < \epsilon.$$

Proof. Let $A = XTX^*$ be the Schur decomposition of A , so that X is a unitary matrix and T is an upper triangular matrix with the eigenvalues of A , $\lambda_1, \dots, \lambda_m$, as its diagonal elements. For any scalar $c > 0$, let the matrix $D_c = \text{diag}(c, c^2, \dots, c^m)$ and note that the diagonal elements of the upper triangular matrix $D_c T D_c^{-1}$ are also $\lambda_1, \dots, \lambda_m$. Further, the i th column sum of $D_c T D_c^{-1}$ is given by

$$\lambda_i + \sum_{j=1}^{i-1} c^{-(i-j)} t_{ji}.$$

Clearly, by choosing c large enough, we can guarantee that

$$\sum_{j=1}^{i-1} |c^{-(i-j)} t_{ij}| < \epsilon$$

for each i . In this case, because $|\lambda_i| \leq \rho(A)$, we must have

$$\|D_c T D_c^{-1}\|_1 < \rho(A) + \epsilon,$$

where $\|A\|_1$ denotes the maximum column sum matrix norm previously defined. For any $m \times m$ matrix B , define $\|B\|_{A,\epsilon}$ as

$$\|B\|_{A,\epsilon} = \|(X D_c^{-1})^{-1} B (X D_c^{-1})\|_1.$$

Since

$$\|A\|_{A,\epsilon} = \|(X D_c^{-1})^{-1} A (X D_c^{-1})\|_1 = \|D_c T D_c^{-1}\|_1,$$

the result follows from Theorem 4.22. \square

Often we will be interested in the limit of a sequence of vectors or the limit of a sequence of matrices. The sequence of $m \times 1$ vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ converges to the $m \times 1$ vector \mathbf{x} if the j th component of \mathbf{x}_k converges to the j th component of \mathbf{x} , as $k \rightarrow \infty$, for each j ; that is, $|x_{jk} - x_j| \rightarrow 0$, as $k \rightarrow \infty$, for each j . Similarly, a sequence of $m \times m$ matrices, A_1, A_2, \dots converges to the $m \times m$ matrix A if each element of A_k converges to the corresponding element of A as $k \rightarrow \infty$. Alternatively, we can consider the notion of the convergence of a sequence with respect to a specific norm. Thus, the sequence of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots$ converges to \mathbf{x} , with respect to the vector norm $\|\mathbf{x}\|$, if $\|\mathbf{x}_k - \mathbf{x}\| \rightarrow 0$ as $k \rightarrow \infty$. Theorem 4.25 indicates that the actual choice of a norm is not important. For a proof of this result, see Horn and Johnson (2013).

Theorem 4.25 Let $\|\mathbf{x}\|_a$ and $\|\mathbf{x}\|_b$ be any two vector norms defined on any $m \times 1$ vector \mathbf{x} . If $\mathbf{x}_1, \mathbf{x}_2, \dots$ is a sequence of $m \times 1$ vectors, then \mathbf{x}_k converges to \mathbf{x} as $k \rightarrow \infty$ with respect to $\|\mathbf{x}\|_a$ if and only if \mathbf{x}_k converges to \mathbf{x} as $k \rightarrow \infty$ with respect to $\|\mathbf{x}\|_b$.

Since the first four conditions of a matrix norm are the conditions of a vector norm, Theorem 4.25 immediately leads to the following.

Corollary 4.25.1 Let $\|A\|_a$ and $\|A\|_b$ be any two matrix norms defined on any $m \times m$ matrix A . If A_1, A_2, \dots is a sequence of $m \times m$ matrices, then A_k converges to A as $k \rightarrow \infty$ with respect to $\|A\|_a$ if and only if A_k converges to A as $k \rightarrow \infty$ with respect to $\|A\|_b$.

A sequence of matrices that is sometimes of interest is the sequence A, A^2, A^3, \dots formed from a fixed $m \times m$ matrix A . A sufficient condition for this sequence of matrices to converge to the null matrix is given next.

Theorem 4.26 Let A be an $m \times m$ matrix, and suppose that for some matrix norm, $\|A\| < 1$. Then $\lim A^k = (0)$ as $k \rightarrow \infty$.

Proof. By repeatedly using condition (e) of a matrix norm, we find that $\|A^k\| \leq \|A\|^k$, and so $\|A^k\| \rightarrow 0$ as $k \rightarrow \infty$, because $\|A\| < 1$. Thus, A^k converges to (0) with respect to the norm $\|A\|$. However, by Corollary 4.25.1, A^k also converges to (0) with respect to the matrix norm (see Problem 4.51)

$$\|A\|_* = m(\max_{1 \leq i, j \leq m} |a_{ij}|).$$

But this implies that $|a_{ij}^k| \rightarrow 0$ as $k \rightarrow \infty$ for each (i, j) , and so the proof is complete. \square

Our next result relates the convergence of A^k to (0) , to the size of the spectral radius of A .

Theorem 4.27 Suppose that A is an $m \times m$ matrix. Then A^k converges to (0) as $k \rightarrow \infty$ if and only if $\rho(A) < 1$.

Proof. Suppose that $A^k \rightarrow (0)$, in which case, $A^k \mathbf{x} \rightarrow \mathbf{0}$ for any $m \times 1$ vector \mathbf{x} . Now if \mathbf{x} is an eigenvector of A corresponding to the eigenvalue λ , we must also have $\lambda^k \mathbf{x} \rightarrow \mathbf{0}$, because $A^k \mathbf{x} = \lambda^k \mathbf{x}$. This can only happen if $|\lambda| < 1$, and so $\rho(A) < 1$, because λ was an arbitrary eigenvalue of A . On the other hand, if $\rho(A) < 1$, then we know from Theorem 4.24 that there is a matrix norm satisfying $\|A\| < 1$. Hence, it follows from Theorem 4.26 that $A^k \rightarrow (0)$. \square

Theorem 4.28 shows that the spectral radius of A is the limit of a particular sequence that can be computed from any matrix norm.

Theorem 4.28 Let A be an $m \times m$ matrix. Then for any matrix norm $\|A\|$

$$\lim_{k \rightarrow \infty} \|A^k\|^{1/k} = \rho(A).$$

Proof. λ is an eigenvalue of A if and only if λ^k is an eigenvalue of A^k . Further, $|\lambda|^k = |\lambda^k|$, so $\rho(A)^k = \rho(A^k)$. This, along with Theorem 4.23, yields $\rho(A)^k \leq \|A^k\|$, or equivalently, $\rho(A) \leq \|A^k\|^{1/k}$. Thus, the proof will be complete if we can show that for arbitrary $\epsilon > 0$, an integer N_ϵ exists, such that $\|A^k\|^{1/k} < \rho(A) + \epsilon$ for all $k > N_\epsilon$. This is the same as showing that an integer N_ϵ exists, such that for all $k > N_\epsilon$, $\|A^k\| < \{\rho(A) + \epsilon\}^k$, or, equivalently,

$$\|B^k\| < 1, \quad (4.14)$$

where $B = \{\rho(A) + \epsilon\}^{-1}A$. However,

$$\rho(B) = \frac{\rho(A)}{\rho(A) + \epsilon} < 1,$$

and so (4.14) follows immediately from Theorem 4.27. \square

Our final result gives a bound that is sometimes useful.

Theorem 4.29 Consider a matrix norm $\|\cdot\|$ and an $m \times m$ matrix A for which $\|I_m\| = 1$ and $\|A\| < 1$. Then $I_m - A$ has an inverse and

$$\|(I_m - A)^{-1}\| \leq \frac{1}{1 - \|A\|}.$$

Proof. Since $\|A\| < 1$, the series $\sum_{i=0}^{\infty} A^i$ converges to some matrix B . Note that

$$(I_m - A) \sum_{i=0}^n A^i = \sum_{i=0}^n A^i - \sum_{i=1}^{n+1} A^i = I_m - A^{n+1}.$$

The limit as $n \rightarrow \infty$ of the left-hand side is $(I_m - A)B$, while the limit of the right-hand side is I_m due to Theorem 4.26. Thus, we have $B = \sum_{i=0}^{\infty} A^i = (I_m - A)^{-1}$, and so

$$\begin{aligned} \|(I_m - A)^{-1}\| &= \|I_m + \sum_{i=1}^{\infty} A^i\| \leq \|I_m\| + \sum_{i=1}^{\infty} \|A^i\| \\ &\leq 1 + \sum_{i=1}^{\infty} \|A\|^i = \frac{1}{1 - \|A\|}, \end{aligned}$$

where the final equality uses the well-known result for a geometric series. \square

Example 4.14 Let A be an $m \times m$ nonsingular matrix. In this example, we wish to compare the inverse of A to the inverse of a perturbation $A + B$ of A , where B is another $m \times m$ matrix satisfying $\|B\| < 1/\|A^{-1}\|$ for some norm for which $\|I_m\| = 1$. Now

$$\|A^{-1}B\| \leq \|A^{-1}\|\|B\| < \|A^{-1}\|/\|A^{-1}\| = 1,$$

and consequently, according to Theorem 4.29, $(I_m + A^{-1}B)$ is nonsingular. Note that

$$\begin{aligned} A^{-1} - (A + B)^{-1} &= \{A^{-1}(A + B) - I_m\}(A + B)^{-1} \\ &= (I_m + A^{-1}B - I_m)(A + B)^{-1} \\ &= A^{-1}B(A + B)^{-1} = A^{-1}B\{A(I_m + A^{-1}B)\}^{-1} \\ &= A^{-1}B(I_m + A^{-1}B)^{-1}A^{-1}. \end{aligned}$$

Using this identity, the bound given in Theorem 4.29, and the fact that $\|A^{-1}B\| \leq \|A^{-1}\|\|B\| < 1$, we then find that a bound on the norm of the error associated with using the inverse of the perturbed matrix instead of the inverse of A is given by

$$\begin{aligned} \|A^{-1} - (A + B)^{-1}\| &\leq \|A^{-1}\|^2\|B\|\|(I_m + A^{-1}B)^{-1}\| \\ &\leq \frac{\|A^{-1}\|^2\|B\|}{1 - \|A^{-1}B\|} \leq \frac{\|A^{-1}\|^2\|B\|}{1 - \|A^{-1}\|\|B\|}. \end{aligned}$$

The corresponding bound on the relative error is

$$\frac{\|A^{-1} - (A + B)^{-1}\|}{\|A^{-1}\|} \leq \frac{\|A^{-1}\|\|B\|}{1 - \|A^{-1}\|\|B\|} = \frac{\kappa(A)\|B\|/\|A\|}{1 - \kappa(A)\|B\|/\|A\|},$$

where $\kappa(A)$, known as the condition number for matrix inversion, is defined as $\kappa(A) = \|A^{-1}\|\|A\|$. Since $\|B\| < 1/\|A^{-1}\|$, it follows that $\kappa(A) < \|A\|/\|B\|$, while a lower bound for $\kappa(A)$ is given by

$$\kappa(A) = \|A\|\|A^{-1}\| \geq \|AA^{-1}\| = \|I_m\| = 1.$$

Clearly, as $\kappa(A)$ increases, the bound on the relative error increases.

PROBLEMS

4.1 Obtain a singular value decomposition for the matrix

$$A = \begin{bmatrix} 1 & 2 & 2 & 1 \\ 1 & 1 & 1 & -1 \end{bmatrix}.$$

4.2 Let A be an $m \times n$ matrix.

- (a) Show that the singular values of A are the same as those of A' .
- (b) Show that the singular values of A are the same as those of FAG , if F and G are orthogonal matrices.
- (c) If $\alpha \neq 0$ is a scalar, how do the singular values of αA compare with those of A ?

4.3 Let A be an $m \times m$ matrix. Show that A has a zero eigenvalue if and only if it has fewer than m singular values.

4.4 Let A be $m \times n$ and B be $n \times m$. We will see in Chapter 7 that the nonzero eigenvalues of AB are the same as those of BA . This is not necessarily true for the singular values. Give an example of matrices A and B for which the singular values of AB are not the same as those of BA .

4.5 Let A be an $m \times n$ matrix having rank r and singular values μ_1, \dots, μ_r . Show that the $(m+n) \times (m+n)$ matrix

$$B = \begin{bmatrix} (0) & A \\ A' & (0) \end{bmatrix}$$

has eigenvalues $\mu_1, \dots, \mu_r, -\mu_1, \dots, -\mu_r$, with the remaining eigenvalues being zero.

4.6 Find a singular value decomposition for the vector $x = (1, 5, 7, 5)'$.

4.7 Let x be an $m \times 1$ nonnull vector and y be an $n \times 1$ nonnull vector. Obtain a singular value decomposition of xy' in terms of x and y .

4.8 Let A be $m \times n$ with $m \leq n$. The polar decomposition of A is

$$A = BR,$$

where B is an $m \times m$ nonnegative definite matrix satisfying $\text{rank}(B) = \text{rank}(A)$ and R is an $m \times n$ matrix satisfying $RR' = I_m$. Use the singular value decomposition to establish the existence of the polar decomposition.

4.9 Let A be an $m \times n$ matrix, and let $A = P_1 \Delta Q_1'$ be the decomposition given in Corollary 4.1.1. Define the $n \times m$ matrix B as $B = Q_1 \Delta^{-1} P_1'$. Simplify, as much as possible, the expressions for ABA and BAB .

4.10 Let A and B be $m \times n$ matrices. From Theorem 1.10 we know that if $B = CAD$, where C and D are nonsingular matrices, then $\text{rank}(B) = \text{rank}(A)$. Prove the converse; that is, if $\text{rank}(B) = \text{rank}(A)$, show that nonsingular matrices C and D exist, such that $B = CAD$.

4.11 If t is an estimator of θ , then the mean squared error (MSE) of t is defined by

$$\text{MSE}(t) = E[(t - \theta)^2] = \text{var}(t) + \{E(t) - \theta\}^2.$$

Consider the multicollinearity problem discussed in Example 4.4 in which r of the singular values of Z_1 are very small relative to the others. Suppose that we

want to estimate the response variable corresponding to an observation that has the standardized explanatory variables at the values given in the $k \times 1$ vector \mathbf{z} . Let $\hat{y} = \bar{y} + \mathbf{z}'(Z_1'Z_1)^{-1}Z_1'\mathbf{y}$ be the estimate obtained using ordinary least squares regression, whereas $\tilde{y} = \bar{y} + \mathbf{z}'U_1D_1^{-1}V_1'\mathbf{y}$ is the estimate obtained using principal components regression, both being estimates of $\theta = \delta_0 + \mathbf{z}'\boldsymbol{\delta}_1$. Assume throughout that $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 I_N)$.

(a) Show that if the vector $\mathbf{v} = (v_1, \dots, v_N)'$ satisfies $\mathbf{z}' = \mathbf{v}'DU'$, then

$$\text{MSE}(\hat{y}) = \sigma^2 \left(N^{-1} + \sum_{i=1}^k v_i^2 \right).$$

(b) Show that

$$\text{MSE}(\tilde{y}) = \sigma^2 \left(N^{-1} + \sum_{i=1}^{k-r} v_i^2 \right) + \left(\sum_{i=k-r+1}^k d_i v_i \alpha_i \right)^2,$$

where d_i is the i th diagonal element of D .

(c) If $r = 1$, when will $\text{MSE}(\tilde{y}) < \text{MSE}(\hat{y})$?

4.12 Suppose that ten observations are obtained in a process involving two explanatory variables and a response variable resulting in the following data:

x_1	x_2	y
-2.49	6.49	28.80
0.85	4.73	21.18
-0.78	4.24	24.73
-0.75	5.54	25.34
1.16	4.74	28.50
-1.52	5.86	27.19
-0.51	5.65	26.22
-0.05	4.50	20.71
-1.01	5.75	25.47
0.13	5.69	29.83

(a) Obtain the matrix of standardized explanatory variables Z_1 , use ordinary least squares to estimate the parameters in the model $\mathbf{y} = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}$, and obtain the fitted values $\hat{\mathbf{y}} = \hat{\delta}_0 \mathbf{1}_N + Z_1 \hat{\boldsymbol{\delta}}_1$.

(b) Compute the spectral decomposition of $Z_1'Z_1$. Then use principal components regression to obtain an alternative vector of fitted values.

- (c) Use both models of (a) and (b) to estimate the response variable for an observation having $x_1 = -2$ and $x_2 = 4$.

4.13 Consider the 3×3 symmetric matrix given by

$$A = \begin{bmatrix} 3 & 1 & -1 \\ 1 & 3 & 1 \\ -1 & 1 & 3 \end{bmatrix}.$$

- (a) Find the spectral decomposition of A .
 (b) Find the symmetric square root matrix for A .
 (c) Find a nonsymmetric square root matrix for A .
- 4.14** Use the spectral decomposition theorem to prove Theorem 4.5
- 4.15** Find a 3×2 matrix T , such that $TT' = A$, where

$$A = \begin{bmatrix} 5 & 4 & 0 \\ 4 & 5 & 3 \\ 0 & 3 & 5 \end{bmatrix}.$$

4.16 Suppose $\mathbf{x} \sim N_3(\mathbf{0}, \Omega)$, where

$$\Omega = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}.$$

Find a 3×3 matrix A , such that the components of $\mathbf{z} = A\mathbf{x}$ are independently distributed.

4.17 Let the matrices A , B , and C be given by

$$A = \begin{bmatrix} 1 & 2 & 5 \\ 2 & 1 & 4 \\ -1 & 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ -1 & 1 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 2 & 1 & -1 \\ 2 & 5 & 3 \\ -2 & -1 & 1 \end{bmatrix}.$$

- (a) Which of these matrices are diagonalizable?
 (b) Which of these matrices have their rank equal to the number of nonzero eigenvalues?
- 4.18** A 3×3 matrix A has eigenvalues $\lambda_1 = 1, \lambda_2 = 2, \lambda_3 = 3$, and corresponding eigenvectors $\mathbf{x}_1 = (1, 1, 2)'$, $\mathbf{x}_2 = (-1, 1, -2)'$ and $\mathbf{x}_3 = (1, 1, 1)'$. Find A .
- 4.19** Let A and B be $m \times m$ matrices and suppose that one of them is nonsingular. Show that if AB is diagonalizable, then BA must also be diagonalizable. Show by example that this is not necessarily true when both A and B are singular.
- 4.20** Let A be an $m \times m$ positive definite matrix and B be an $m \times m$ symmetric matrix. Show that AB is a diagonalizable matrix and the number of its positive, negative, and zero eigenvalues are the same as that of B .

4.21 Let A be an $m \times m$ matrix and B be an $n \times n$ matrix. Prove that the matrix

$$C = \begin{bmatrix} A & (0) \\ (0) & B \end{bmatrix}$$

is diagonalizable if and only if the matrices A and B are diagonalizable. Using induction, show that the square matrices A_1, \dots, A_k are diagonalizable if and only if $\text{diag}(A_1, \dots, A_k)$ is diagonalizable.

4.22 Let \mathbf{x} and \mathbf{y} be $m \times 1$ vectors. Show that $A = \mathbf{x}\mathbf{y}'$ is diagonalizable if and only if $\mathbf{x}'\mathbf{y} \neq 0$.

4.23 Let $A = \sum_{i=1}^k \mu_i P_A(\mu_i)$ be the spectral decomposition of the diagonalizable matrix A . Show that

(a) $P_A(\mu_i)P_A(\mu_j) = (0)$, if $i \neq j$,

(b) $\sum_{i=1}^k P_A(\mu_i) = I_m$,

(c) $P_A(\mu_i)$ is the projection matrix for the projection onto the null space of $(A - \mu_i I_m)$ along the column space of $(A - \mu_i I_m)$.

4.24 Suppose A is an $m \times m$ diagonalizable matrix with linearly independent eigenvectors given by the columns of the $m \times m$ matrix X . Show that linearly independent eigenvectors of A' are given by the columns of $Y = X^{-1'}$.

4.25 Find a 4×4 matrix A having eigenvalues 0 and 1 with multiplicities 3 and 1, respectively, such that

(a) the rank of A is 1,

(b) the rank of A is 2,

(c) the rank of A is 3.

4.26 Repeat Example 4.10 for 5×5 matrices; that is, obtain a collection of 5×5 matrices in Jordan canonical form, such that every 5×5 matrix having the eigenvalue λ with multiplicity 5 is similar to one of the matrices in this set.

4.27 Consider the 6×6 matrix

$$J = \begin{bmatrix} 2 & 1 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 1 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 3 & 1 \\ 0 & 0 & 0 & 0 & 0 & 3 \end{bmatrix},$$

which is in Jordan canonical form.

(a) Find the eigenvalues of J and their multiplicities.

(b) Find the eigenspaces of J .

4.28 An $m \times m$ matrix B is said to be nilpotent if $B^k = (0)$ for some positive integer k .

(a) Show that $J_h(\lambda) = \lambda I_h + B_h$, where B_h is nilpotent. In particular, show that $B_h^h = (0)$.

- (b) Let $J = \text{diag}(J_{h_1}(\lambda_1), \dots, J_{h_r}(\lambda_r))$ be a Jordan canonical form. Show that J can be written as $J = D + B$, where D is diagonal and B is nilpotent. What is the smallest h such that $B^h = (0)$?
- (c) Use part (b) to show that if A is similar to J , then A can be expressed as $A = F + G$, where F is diagonalizable and G is nilpotent.
- 4.29** Suppose that A is 5×5 with the eigenvalue λ having multiplicity 5. If $(A - \lambda I_5)^2 = (0)$, what are the possible Jordan canonical forms for A ?
- 4.30** If $J_h(\lambda)$ is an $h \times h$ Jordan block matrix, find the eigenvalues of $\{J_h(\lambda)\}^2$. If $\lambda \neq 0$, how many linearly independent eigenvectors does $\{J_h(\lambda)\}^2$ have? Use this information to show that if an $m \times m$ nonsingular matrix A has Jordan decomposition $A = B^{-1}JB$, where $J = \text{diag}(J_{h_1}(\lambda_1), \dots, J_{h_r}(\lambda_r))$, then A^2 has the Jordan decomposition given by $A^2 = B_*^{-1}J_*B_*$, where $J_* = \text{diag}(J_{h_1}(\lambda_1^2), \dots, J_{h_r}(\lambda_r^2))$.
- 4.31** Let A be an $m \times m$ nilpotent matrix. In Problem 3.39, it was shown that all of the eigenvalues of A are 0. Use this and the Jordan canonical form of A to show that there must be a positive integer $h \leq m$ satisfying $A^h = (0)$.
- 4.32** Let A be an $m \times m$ matrix. Show that the rank of A is equal to the number of nonzero eigenvalues of A if and only if $\text{rank}(A^2) = \text{rank}(A)$.
- 4.33** Suppose that λ is an eigenvalue of A with multiplicity r . Show that there are r linearly independent eigenvectors of A corresponding to λ if and only if $\text{rank}(A - \lambda I_m) = \text{rank}\{(A - \lambda I_m)^2\}$.
- 4.34** Let A and B be $m \times m$ matrices. Suppose that an $m \times m$ unitary matrix X exists, such X^*AX and X^*BX are both upper triangular matrices. Show then that the eigenvalues of $AB - BA$ are all equal to 0.
- 4.35** Let T and U be $m \times m$ upper triangular matrices. In addition, suppose that for some positive integer $r < m$, $t_{ij} = 0$ for $1 \leq i \leq r$, $1 \leq j \leq r$, and $u_{r+1, r+1} = 0$. Show that the upper triangular matrix $V = TU$ is such that $v_{ij} = 0$ for $1 \leq i \leq r + 1$, $1 \leq j \leq r + 1$.
- 4.36** Use the Schur decomposition of a matrix A and the result of the previous exercise to prove the Cayley–Hamilton theorem given as Theorem 3.8; that is, if $\lambda_1, \dots, \lambda_m$ are the eigenvalues of A , show that

$$(A - \lambda_1 I_m)(A - \lambda_2 I_m) \cdots (A - \lambda_m I_m) = (0).$$

- 4.37** Obtain a Schur decomposition for the matrix C given in Problem 4.17.
- 4.38** Repeat Problem 4.37 by obtaining a different Schur decomposition of C .
- 4.39** Let A be an $m \times m$ matrix, and consider the Schur decomposition given in Theorem 4.12. Show that although the matrix T is not uniquely defined, the quantity $\sum_{i < j} |t_{ij}|^2$ is uniquely defined.
- 4.40** Let A and B be $m \times m$ matrices and suppose that $AB = BA$. Show that there exists an $m \times m$ unitary matrix X such that X^*AX and X^*BX are upper triangular.

4.41 Suppose that A and B are $m \times m$ and diagonalizable. Show that A and B commute, that is, $AB = BA$, if and only if they are simultaneously diagonalizable; in other words, $AB = BA$, if and only if a nonsingular matrix X exists, such that both $X^{-1}AX$ and $X^{-1}BX$ are diagonal matrices. This proves Theorem 4.20 when $k = 2$.

4.42 Let

$$A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}.$$

(a) Show that $AB = BA$.

(b) Show that AB is not diagonalizable.

(c) Why does this not contradict the result of Problem 4.41?

4.43 Suppose that the $m \times m$ matrices A and B are diagonalizable and $AB = BA$. Denote the eigenvalues of A by $\lambda_1, \dots, \lambda_m$ and those of B by μ_1, \dots, μ_m . If the eigenvalues of $A + B$ are $\gamma_1, \dots, \gamma_m$, show that for $k = 1, \dots, m$,

$$\gamma_k = \lambda_{i_k} + \mu_{j_k},$$

where (i_1, \dots, i_m) and (j_1, \dots, j_m) are permutations of $(1, \dots, m)$.

4.44 Let A and B be $m \times m$ matrices, and suppose that A and B commute.

(a) If A and B are nonsingular, show that A^{-1} and B^{-1} commute.

(b) If i and j are positive integers, show that A^i and B^j commute.

4.45 Find 2×2 matrices A and B that do not satisfy the conditions of Theorem 4.15 and Theorem 4.16, yet a nonsingular matrix C exists for which CAC' and CBC' are diagonal.

4.46 Suppose that A and B are $m \times m$ positive definite matrices. Show that $A - B$ is positive definite if and only if $B^{-1} - A^{-1}$ is positive definite.

4.47 Let A and B be $m \times m$ symmetric matrices with B being positive definite. Show that $A + B$ is positive definite if and only if every eigenvalue of AB^{-1} is greater than -1 .

4.48 Let A and B be $m \times m$ matrices, and suppose A has m distinct eigenvalues. Show that if $AB = BA$, then there exists a nonsingular matrix X such that both $X^{-1}AX$ and $X^{-1}BX$ are diagonal.

4.49 Suppose A and B are $m \times m$ symmetric matrices.

(a) Show that if B is positive definite, then AB is diagonalizable.

(b) Show that if both A and B are nonnegative definite, then AB is diagonalizable.

4.50 Show that the functions $\|A\|_E$, $\|A\|_1$, $\|A\|_\infty$, and $\|A\|_2$ given in Section 4.8 are, in fact, matrix norms.

4.51 Let A be an $m \times m$ matrix, and consider the function

$$\|A\|_* = m \left(\max_{1 \leq i, j \leq m} |a_{ij}| \right).$$

Show that $\|A\|_*$ is a matrix norm.

4.52 Prove Theorem 4.22.

4.53 For any matrix norm defined on $m \times m$ matrices, show that

(a) $\|I_m\| \geq 1$,

(b) $\|A^{-1}\| \geq \|A\|^{-1}$, if A is an $m \times m$ nonsingular matrix.

4.54 Show that

(a) the maximum column sum matrix norm, $\|A\|_1$, is induced by the sum vector norm, $\|x\|_1$,

(b) the maximum row sum matrix norm, $\|A\|_\infty$, is induced by the infinity vector norm, $\|x\|_\infty$.

4.55 Let A be an $m \times m$ real matrix with singular values $\delta_1, \dots, \delta_r$. Show that

$$\|A\|_E = \left(\sum_{i=1}^r \delta_i^2 \right)^{1/2}.$$

4.56 Let A be an $m \times m$ real matrix with singular values $\delta_1 \geq \dots \geq \delta_r$, and suppose that B is another $m \times m$ real matrix with $\text{rank}(B) = s < r$. Show that

$$\|B - A\|_E^2 \geq \sum_{i=s+1}^r \delta_i^2.$$

4.57 Find examples of 2×2 matrices A and B to show that it is possible to have

(a) $\rho(A + B) > \rho(A) + \rho(B)$,

(b) $\rho(AB) > \rho(A)\rho(B)$.

4.58 Consider the 2×2 matrix of the form

$$A = \begin{bmatrix} a & 1 \\ 0 & a \end{bmatrix}.$$

(a) Determine A^k for general positive integer k .

(b) Find $\rho(A)$ and $\rho(A^k)$.

(c) For which values of a does A^k converge to (0) as $k \rightarrow \infty$? In this case, show how to construct a norm so that $\|A\| < 1$.

4.59 In this problem, we consider a factorization of an $m \times m$ matrix A of the form $A = LU$, where L is an $m \times m$ lower triangular matrix and U is an $m \times m$ upper triangular matrix.

(a) Let A_j be the $j \times j$ submatrix of A consisting of the first j rows and j columns of A . Show that if $r = \text{rank}(A)$ and $|A_j| \neq 0, j = 1, \dots, r$, then A_r can be factored as $A_r = L_* U_*$, where L_* is an $r \times r$ nonsingular lower triangular matrix and U_* is an $r \times r$ nonsingular upper triangular matrix. Apply this result to then show that A may be factored as $A = LU$, where L

is an $m \times m$ lower triangular matrix and U is an $m \times m$ upper triangular matrix.

- (b) Show that not every $m \times m$ matrix has an LU factorization by finding a 2×2 matrix that cannot be factored in this way.
 - (c) Show how the LU factorization of A can be used to simplify the computation of a solution \mathbf{x} , to the system of equations $A\mathbf{x} = \mathbf{c}$.
- 4.60** Suppose that A is an $m \times m$ matrix. Show that an $m \times m$ lower triangular matrix L , an $m \times m$ upper triangular matrix U , and $m \times m$ permutation matrices P and Q exist, such that $A = PLUQ$.
- 4.61** Suppose that A is an $m \times m$ matrix for which $|A_j| \neq 0, j = 1, \dots, m$, where A_j denotes the $j \times j$ submatrix of A consisting of the first j rows and j columns of A .
- (a) Show that there exist $m \times m$ lower triangular matrices L and M having all diagonal elements equal to one and an $m \times m$ diagonal matrix D , such that $A = LDM'$.
 - (b) Show that if A is also symmetric, then $M = L$, so that $A = LDL'$.

5

GENERALIZED INVERSES

5.1 INTRODUCTION

The inverse of a matrix is defined for all square matrices that are nonsingular. There are some situations in which we may have a rectangular matrix or a square singular matrix, A , and still be in need of another matrix that in some ways behaves like the inverse of A . One such situation, which is often encountered in the study of statistics as well as in many other fields of application, involves finding solutions to a system of linear equations. A system of linear equations can be written in matrix form as

$$A\mathbf{x} = \mathbf{c},$$

where A is an $m \times n$ matrix of constants, \mathbf{c} is an $m \times 1$ vector of constants, and \mathbf{x} is an $n \times 1$ vector of variables for which we need to find solutions. If $m = n$ and A is nonsingular, then A^{-1} exists, and so by premultiplying our system of equations by A^{-1} , we see that the system is satisfied only if $\mathbf{x} = A^{-1}\mathbf{c}$; that is, the system has a solution, the solution is unique, and it is given by $\mathbf{x} = A^{-1}\mathbf{c}$. When A^{-1} does not exist, how do we determine whether the system has any solutions, and if solutions exist, how many solutions are there, and how do we find them? We will see in Chapter 6 that the answers to all of these questions can be conveniently expressed in terms of the generalized inverses discussed in this chapter.

A second application of generalized inverses in statistics involves quadratic forms and chi-squared distributions. Suppose we have an m -dimensional random vector \mathbf{x}

that has a mean vector of zero and covariance matrix Ω . A useful transformation in some situations is one that transforms \mathbf{x} to another random vector, \mathbf{z} , having the identity matrix as its covariance matrix. For instance, in Chapter 11, we will see that if \mathbf{z} has a normal distribution, then the sum of squares of the components of \mathbf{z} , that is, $\mathbf{z}'\mathbf{z}$, has a chi-squared distribution. We saw in Example 4.5 that if T is any $m \times m$ matrix satisfying $\Omega^{-1} = TT'$, then $\mathbf{z} = T'\mathbf{x}$ will have I_m as its covariance matrix. Then

$$\mathbf{z}'\mathbf{z} = \mathbf{x}'(T')'T'\mathbf{x} = \mathbf{x}'(TT')\mathbf{x} = \mathbf{x}'\Omega^{-1}\mathbf{x},$$

which of course, will be possible only if Ω is positive definite. If Ω is positive semidefinite with rank r , then it will be possible to find $m \times m$ matrices A and B , with $A = BB'$, such that when \mathbf{z} is defined by $\mathbf{z} = B'\mathbf{x}$,

$$\text{var}(\mathbf{z}) = \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix},$$

and $\mathbf{z}'\mathbf{z} = \mathbf{x}'A\mathbf{x}$. We will see later that A is a generalized inverse of Ω and $\mathbf{z}'\mathbf{z}$ still has a chi-squared distribution if \mathbf{z} has a normal distribution.

5.2 THE MOORE–PENROSE GENERALIZED INVERSE

A useful generalized inverse in statistical applications is one developed by Moore (1920, 1935) and Penrose (1955). This inverse is defined so as to possess four properties that the inverse of a square nonsingular matrix has.

Definition 5.1 The Moore–Penrose inverse of the $m \times n$ matrix A is the $n \times m$ matrix, denoted by A^+ , which satisfies the conditions

$$AA^+A = A, \tag{5.1}$$

$$A^+AA^+ = A^+, \tag{5.2}$$

$$(AA^+)' = AA^+, \tag{5.3}$$

$$(A^+A)' = A^+A. \tag{5.4}$$

One of the most important features of the Moore–Penrose inverse, one that distinguishes it from other generalized inverses that we will discuss in this chapter, is that it is uniquely defined. This fact, along with the existence of the Moore–Penrose inverse, is established in Theorem 5.1.

Theorem 5.1 Corresponding to each $m \times n$ matrix A , one and only one $n \times m$ matrix A^+ exists satisfying conditions (5.1)–(5.4).

Proof. First we will prove the existence of A^+ . If A is the $m \times n$ null matrix, then it is easily verified that the four conditions in Definition 5.1 are satisfied with

$A^+ = (0)$, the $n \times m$ null matrix. If $A \neq (0)$, so that $\text{rank}(A) = r > 0$, then from Corollary 4.1.1, we know $m \times r$ and $n \times r$ matrices P and Q exist, such that $P'P = Q'Q = I_r$ and

$$A = P\Delta Q',$$

where Δ is a diagonal matrix with positive diagonal elements. Note that if we define $A^+ = Q\Delta^{-1}P'$, then

$$\begin{aligned} AA^+A &= P\Delta Q'Q\Delta^{-1}P'P\Delta Q' = P\Delta\Delta^{-1}\Delta Q' \\ &= P\Delta Q' = A, \\ A^+AA^+ &= Q\Delta^{-1}P'P\Delta Q'Q\Delta^{-1}P' = Q\Delta^{-1}\Delta\Delta^{-1}P' \\ &= Q\Delta^{-1}P' = A^+, \\ AA^+ &= P\Delta Q'Q\Delta^{-1}P' = PP' \quad \text{is symmetric,} \\ A^+A &= Q\Delta^{-1}P'P\Delta Q' = QQ' \quad \text{is symmetric.} \end{aligned}$$

Thus, $A^+ = Q\Delta^{-1}P'$ is a Moore–Penrose inverse of A , and so we have established the existence of the Moore–Penrose inverse. Next, suppose that B and C are any two matrices satisfying conditions (5.1)–(5.4) for A^+ . Then using these four conditions, we find that

$$\begin{aligned} AB &= (AB)' = B'A' = B'(ACA)' = B'A'(AC)' \\ &= (AB)'AC = ABAC = AC \end{aligned}$$

and

$$\begin{aligned} BA &= (BA)' = A'B' = (ACA)'B' = (CA)'A'B' \\ &= CA(BA)' = CABA = CA. \end{aligned}$$

Now using these two identities and (5.2), we see that

$$B = BAB = BAC = CAC = C.$$

Since B and C are identical, the Moore–Penrose inverse is unique. \square

We saw in the proof of Theorem 5.1 that the Moore–Penrose inverse of a matrix A is explicitly related to the singular value decomposition of A ; that is, this inverse is nothing more than a simple function of the component matrices making up the singular value decomposition of A .

Definition 5.1 is the definition of a generalized inverse given by Penrose (1955). The following alternative definition, which we will find useful on some occasions, utilizes properties of the Moore–Penrose inverse that were first illustrated by

Moore (1935). This definition applies the concept of projection matrices that was discussed in Chapter 2. Recall that if S is a vector subspace of R^m and P_S is its projection matrix, then for any $\mathbf{x} \in R^m$, $P_S \mathbf{x}$ gives the orthogonal projection of \mathbf{x} onto S , whereas $\mathbf{x} - P_S \mathbf{x}$ is the component of \mathbf{x} orthogonal to S ; further, the unique matrix P_S is given by $\mathbf{x}_1 \mathbf{x}'_1 + \cdots + \mathbf{x}_r \mathbf{x}'_r$, where $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ is any orthonormal basis for S .

Definition 5.2 Let A be an $m \times n$ matrix. Then the Moore–Penrose inverse of A is the unique $n \times m$ matrix A^+ satisfying

- (a) $AA^+ = P_{R(A)}$,
- (b) $A^+A = P_{R(A^+)}$,

where $P_{R(A)}$ and $P_{R(A^+)}$ are the projection matrices of the range spaces of A and A^+ , respectively.

It should be noted that although $A^+AA' = A'$, condition (b) in Definition 5.2 cannot simply be replaced by the condition $A^+A = P_{R(A')}$. As shown by Hu (2008), this substitution would then require a third condition, $\text{rank}(A^+) = \text{rank}(A)$.

The equivalence of Definition 5.1 and Definition 5.2 is not immediately obvious. Consequently, we will establish it in Theorem 5.2.

Theorem 5.2 Definition 5.2 is equivalent to Definition 5.1.

Proof. We first show that a matrix A^+ satisfying Definition 5.2 must also satisfy Definition 5.1. Conditions (5.3) and (5.4) follow immediately because by definition, a projection matrix is symmetric, whereas (5.1) and (5.2) follow because the columns of A are in $R(A)$ imply that

$$AA^+A = P_{R(A)}A = A,$$

and the columns of A^+ are in $R(A^+)$ imply that

$$A^+AA^+ = P_{R(A^+)}A^+ = A^+.$$

Conversely, now suppose that A^+ satisfies Definition 5.1. Premultiplying (5.2) by A yields the identity

$$AA^+AA^+ = (AA^+)^2 = AA^+,$$

which along with (5.3) shows that AA^+ is idempotent and symmetric, and thus by Theorem 2.22 is a projection matrix. To show that it is the projection matrix of the range space of A , note that for any matrices B and C , for which BC is defined, $R(BC) \subseteq R(B)$. Using this twice along with (5.1), we find that

$$R(A) = R(AA^+A) \subseteq R(AA^+) \subseteq R(A),$$

so that $R(AA^+) = R(A)$. This proves that $P_{R(A)} = AA^+$. A proof of $P_{R(A^+)} = A^+A$ is obtained in a similar fashion using (5.1) and (5.4). \square

5.3 SOME BASIC PROPERTIES OF THE MOORE–PENROSE INVERSE

In this section, we will establish some of the basic properties of the Moore–Penrose inverse, whereas in some of the subsequent sections we will look at some more specialized results. First, we have Theorem 5.3.

Theorem 5.3 Let A be an $m \times n$ matrix. Then

- (a) $(\alpha A)^+ = \alpha^{-1}A^+$, if $\alpha \neq 0$ is a scalar,
- (b) $(A')^+ = (A^+)',$
- (c) $(A^+)^+ = A,$
- (d) $A^+ = A^{-1}$, if A is square and nonsingular,
- (e) $(A'A)^+ = A^+A^{+'}$ and $(AA')^+ = A^{+'}A^+,$
- (f) $(AA^+)^+ = AA^+$ and $(A^+A)^+ = A^+A,$
- (g) $A^+ = (A'A)^+A' = A'(AA^+)^+,$
- (h) $A^+ = (A'A)^{-1}A'$ and $A^+A = I_n$, if $\text{rank}(A) = n,$
- (i) $A^+ = A'(AA')^{-1}$ and $AA^+ = I_m$, if $\text{rank}(A) = m,$
- (j) $A^+ = A'$, if the columns of A are orthogonal, that is, $A'A = I_n.$

Proof. Each part is proven by simply verifying that the stated inverse satisfies conditions (5.1)–(5.4). Here, we will only verify that $(A'A)^+ = A^+A^{+'}$, given in (e), and leave the remaining proofs to the reader. Since A^+ satisfies the four conditions of a Moore–Penrose inverse, we find that

$$\begin{aligned}
 A'A(A'A)^+A'A &= A'AA^+A^{+'}A'A = A'AA^+(AA^+)'A \\
 &= A'AA^+AA^+A = A'AA^+A = A'A, \\
 (A'A)^+A'A(A'A)^+ &= A^+A^{+'}A'AA^+A^{+'} = A^+(AA^+)'AA^+A^{+'} \\
 &= A^+AA^+AA^+A^{+'} = A^+AA^+A^{+'} \\
 &= A^+A^{+'} = (A'A)^+,
 \end{aligned}$$

so that $A^+A^{+'}$ satisfies conditions (5.1) and (5.2) of the Moore–Penrose inverse $(A'A)^+$. In addition, note that

$$\begin{aligned}
 A'A(A'A)^+ &= A'AA^+A^{+'} = A'(A^+(AA^+))' \\
 &= A'(A^+AA^+)' = A'A^{+'} \\
 &= (A^+A)'.
 \end{aligned}$$

and A^+A must be symmetric by definition, so it follows that condition (5.3) is satisfied for $(A'A)^+ = A^+A^{+'}$. Likewise, condition (5.4) holds because

$$\begin{aligned}(A'A)^+A'A &= A^+A^{+'}A'A = A^+(AA^+)'A \\ &= A^+AA^+A = A^+A.\end{aligned}$$

This then proves that $(A'A)^+ = A^+A^{+'}$. □

Example 5.1 Properties (h) and (i) of Theorem 5.3 give useful ways of computing the Moore–Penrose inverse of matrices that have full column rank or full row rank. We will demonstrate this by finding the Moore–Penrose inverses of

$$\mathbf{a} = \begin{bmatrix} 1 \\ 1 \end{bmatrix} \quad \text{and} \quad A = \begin{bmatrix} 1 & 2 & 1 \\ 2 & 1 & 0 \end{bmatrix}.$$

From property (h), for any vector $\mathbf{a} \neq \mathbf{0}$, \mathbf{a}^+ will be given by $(\mathbf{a}'\mathbf{a})^{-1}\mathbf{a}'$, so here we find that

$$\mathbf{a}^+ = [0.5 \quad 0.5].$$

For A , we can use property (i) because $\text{rank}(A) = 2$. Computing AA' and $(AA')^{-1}$, we get

$$AA' = \begin{bmatrix} 6 & 4 \\ 4 & 5 \end{bmatrix}, \quad (AA')^{-1} = \frac{1}{14} \begin{bmatrix} 5 & -4 \\ -4 & 6 \end{bmatrix},$$

and so

$$\begin{aligned}A^+ &= A'(AA')^{-1} = \frac{1}{14} \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 5 & -4 \\ -4 & 6 \end{bmatrix} \\ &= \frac{1}{14} \begin{bmatrix} -3 & 8 \\ 6 & -2 \\ 5 & -4 \end{bmatrix}.\end{aligned}$$

Theorem 5.4 establishes a relationship between the rank of a matrix and the rank of its Moore–Penrose inverse.

Theorem 5.4 For any $m \times n$ matrix A ,

$$\text{rank}(A) = \text{rank}(A^+) = \text{rank}(AA^+) = \text{rank}(A^+A).$$

Proof. Using condition (5.1) and the fact that the rank of a matrix product cannot exceed the rank of any of the matrices in the product, we find that

$$\text{rank}(A) = \text{rank}(AA^+A) \leq \text{rank}(AA^+) \leq \text{rank}(A^+). \quad (5.5)$$

In a similar fashion, using condition (5.2), we get

$$\text{rank}(A^+) = \text{rank}(A^+AA^+) \leq \text{rank}(A^+A) \leq \text{rank}(A). \quad (5.6)$$

The result follows immediately from (5.5) and (5.6). \square

We have seen through Definition 5.2 and Theorem 5.2 that A^+A is the projection matrix of the range of A^+ . It also will be the projection matrix of the range of any matrix B satisfying $\text{rank}(B) = \text{rank}(A^+)$ and $A^+AB = B$. For instance, from Theorem 5.4, we have $\text{rank}(A') = \text{rank}(A^+)$ and

$$A^+AA' = (A^+A)'A' = A'A^+A' = A',$$

so A^+A is also the projection matrix of the range of A' ; that is, $P_{R(A')} = A^+A$.

Theorem 5.5 summarizes some of the special properties possessed by the Moore–Penrose inverse of a symmetric matrix.

Theorem 5.5 Let A be an $m \times m$ symmetric matrix. Then

- (a) A^+ is also symmetric,
- (b) $AA^+ = A^+A$,
- (c) $A^+ = A$, if A is idempotent.

Proof. Using Theorem 5.3(b) and the fact that $A = A'$, we have

$$A^+ = (A')^+ = (A^+)',$$

which then proves (a). To prove (b), note that it follows from condition (5.3) of the Moore–Penrose inverse of a matrix, along with the symmetry of both A and A^+ , that

$$AA^+ = (AA^+)' = A^{+'}A' = A^+A.$$

Finally, (c) is established by verifying the four conditions of the Moore–Penrose inverse for $A^+ = A$, when $A^2 = A$. For instance, both conditions (5.1) and (5.2) hold because

$$AAA = A^2A = AA = A^2 = A.$$

Note also that

$$(AA)' = A'A' = AA,$$

so that conditions (5.3) and (5.4) hold as well. \square

In the proof of Theorem 5.1, we saw that the Moore–Penrose inverse of any matrix can be conveniently expressed in terms of the components involved in the singular

value decomposition of that matrix. Likewise, in the special case of a symmetric matrix, we will be able to write the Moore–Penrose inverse in terms of the components of the spectral decomposition of that matrix, that is, in terms of its eigenvalues and eigenvectors. Before identifying this relationship, we first consider the Moore–Penrose inverse of a diagonal matrix. The proof of this result, which simply involves the verification of conditions (5.1)–(5.4), is left to the reader.

Theorem 5.6 Let Λ be the $m \times m$ diagonal matrix $\text{diag}(\lambda_1, \dots, \lambda_m)$. Then the Moore–Penrose inverse Λ^+ of Λ is the diagonal matrix $\text{diag}(\phi_1, \dots, \phi_m)$, where

$$\phi_i = \begin{cases} \lambda_i^{-1}, & \text{if } \lambda_i \neq 0, \\ 0, & \text{if } \lambda_i = 0. \end{cases}$$

Theorem 5.7 Let $\mathbf{x}_1, \dots, \mathbf{x}_m$ be a set of orthonormal eigenvectors corresponding to the eigenvalues, $\lambda_1, \dots, \lambda_m$, of the $m \times m$ symmetric matrix A . If we define $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ and $X = (\mathbf{x}_1, \dots, \mathbf{x}_m)$, then

$$A^+ = X\Lambda^+X'.$$

Proof. Let $r = \text{rank}(A)$, and suppose that we have ordered the λ_i 's so that $\lambda_{r+1} = \dots = \lambda_m = 0$. Partition X as $X = [X_1 \ X_2]$, where X_1 is $m \times r$, and partition Λ in block diagonal form as $\Lambda = \text{diag}(\Lambda_1, (0))$, where $\Lambda_1 = \text{diag}(\lambda_1, \dots, \lambda_r)$. Then, the spectral decomposition of A is given by

$$A = [X_1 \ X_2] \begin{bmatrix} \Lambda_1 & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} X_1' \\ X_2' \end{bmatrix} = X_1\Lambda_1X_1',$$

and similarly the expression above for A^+ reduces to $A^+ = X_1\Lambda_1^{-1}X_1'$. Thus, because $X_1'X_1 = I_r$, we have

$$AA^+ = X_1\Lambda_1X_1'X_1\Lambda_1^{-1}X_1' = X_1\Lambda_1\Lambda_1^{-1}X_1' = X_1X_1',$$

which is clearly symmetric, so condition (5.3) is satisfied. Similarly, $A^+A = X_1X_1'$, and so (5.4) also holds. Conditions (5.1) and (5.2) hold because

$$\begin{aligned} AA^+A &= (AA^+)A = X_1X_1'X_1\Lambda_1X_1' \\ &= X_1\Lambda_1X_1' = A \end{aligned}$$

and

$$\begin{aligned} A^+AA^+ &= A^+(AA^+) = X_1\Lambda_1^{-1}X_1'X_1X_1' \\ &= X_1\Lambda_1^{-1}X_1' = A^+, \end{aligned}$$

and so the proof is complete. \square

Example 5.2 Consider the symmetric matrix

$$A = \begin{bmatrix} 32 & 16 & 16 \\ 16 & 14 & 2 \\ 16 & 2 & 14 \end{bmatrix}.$$

It is easily verified that an eigenanalysis of A reveals that it can be expressed as

$$A = \begin{bmatrix} 2/\sqrt{6} & 0 \\ 1/\sqrt{6} & -1/\sqrt{2} \\ 1/\sqrt{6} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 48 & 0 \\ 0 & 12 \end{bmatrix} \begin{bmatrix} 2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix}.$$

Thus, using Theorem 5.7, we find that

$$\begin{aligned} A^+ &= \begin{bmatrix} 2/\sqrt{6} & 0 \\ 1/\sqrt{6} & -1/\sqrt{2} \\ 1/\sqrt{6} & 1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 1/48 & 0 \\ 0 & 1/12 \end{bmatrix} \begin{bmatrix} 2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & -1/\sqrt{2} & 1/\sqrt{2} \end{bmatrix} \\ &= \frac{1}{288} \begin{bmatrix} 4 & 2 & 2 \\ 2 & 13 & -11 \\ 2 & -11 & 13 \end{bmatrix}. \end{aligned}$$

In Section 2.7, we saw that if the columns of an $m \times r$ matrix X form a basis for a vector space S , then the projection matrix of S is given by $X(X'X)^{-1}X'$; that is

$$P_{R(X)} = X(X'X)^{-1}X'.$$

Definition 5.2 indicates how this can be generalized to the situation in which X is not full column rank. Thus, using Definition 5.2 and Theorem 5.3(g), we find that the projection matrix of the space spanned by the columns of X is

$$P_{R(X)} = XX^+ = X(X'X)^+X'. \quad (5.7)$$

Example 5.3 We will use (5.7) to obtain the projection matrix of the range of

$$X = \begin{bmatrix} 4 & 1 & 3 \\ -4 & -3 & -1 \\ 0 & -2 & 2 \end{bmatrix}.$$

The Moore–Penrose inverse of

$$X'X = \begin{bmatrix} 32 & 16 & 16 \\ 16 & 14 & 2 \\ 16 & 2 & 14 \end{bmatrix}$$

was obtained in the previous example, which we use to find that

$$\begin{aligned}
 P_{R(X)} &= X(X'X)^+X' \\
 &= \frac{1}{288} \begin{bmatrix} 4 & 1 & 3 \\ -4 & -3 & -1 \\ 0 & -2 & 2 \end{bmatrix} \begin{bmatrix} 4 & 2 & 2 \\ 2 & 13 & -11 \\ 2 & -11 & 13 \end{bmatrix} \begin{bmatrix} 4 & -4 & 0 \\ 1 & -3 & -2 \\ 3 & -1 & 2 \end{bmatrix} \\
 &= \frac{1}{3} \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}.
 \end{aligned}$$

This illustrates the use of (5.7). Actually, $P_{R(X)}$ can be computed without ever formally computing any Moore–Penrose inverse because $P_{R(X)}$ is the total eigenprojection corresponding to the positive eigenvalues of XX' . Here we have

$$XX' = \begin{bmatrix} 26 & -22 & 4 \\ -22 & 26 & 4 \\ 4 & 4 & 8 \end{bmatrix},$$

which has $\mathbf{z}_1 = (1/\sqrt{2}, -1/\sqrt{2}, 0)'$ and $\mathbf{z}_2 = (1/\sqrt{6}, 1/\sqrt{6}, 2/\sqrt{6})'$ as normalized eigenvectors corresponding to its two positive eigenvalues. Thus, if we let $Z = (\mathbf{z}_1, \mathbf{z}_2)$, then

$$P_{R(X)} = ZZ' = \frac{1}{3} \begin{bmatrix} 2 & -1 & 1 \\ -1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}.$$

Example 5.4 The Moore–Penrose inverse is useful in constructing quadratic forms in normal random vectors so that they have chi-squared distributions. This is a topic that we will investigate in more detail in Chapter 11; here we will look at a simple illustration. A common situation encountered in inferential statistics is one in which one has a sample statistic, $\mathbf{t} \sim N_m(\boldsymbol{\theta}, \Omega)$, and it is desired to determine whether the $m \times 1$ parameter vector $\boldsymbol{\theta} = \mathbf{0}$; formally, we want to test the null hypothesis $H_0: \boldsymbol{\theta} = \mathbf{0}$ versus the alternative hypothesis $H_1: \boldsymbol{\theta} \neq \mathbf{0}$. One approach to this problem, if Ω is positive definite, is to base the decision between H_0 and H_1 on the statistic

$$v_1 = \mathbf{t}'\Omega^{-1}\mathbf{t}.$$

Now if T is any $m \times m$ matrix satisfying $TT' = \Omega$, and we define $\mathbf{u} = T^{-1}\mathbf{t}$, then $E(\mathbf{u}) = T^{-1}\boldsymbol{\theta}$ and

$$\text{var}(\mathbf{u}) = T^{-1}\{\text{var}(\mathbf{t})\}T^{-1'} = T^{-1}(TT')T^{-1'} = I_m,$$

so $\mathbf{u} \sim N_m(T^{-1}\boldsymbol{\theta}, I_m)$. Consequently, u_1, \dots, u_m are independently distributed normal random variables, and so

$$v_1 = \mathbf{t}'\Omega^{-1}\mathbf{t} = \mathbf{u}'\mathbf{u} = \sum_{i=1}^m u_i^2$$

has a chi-squared distribution with m degrees of freedom. This chi-squared distribution is central if $\boldsymbol{\theta} = \mathbf{0}$ and noncentral if $\boldsymbol{\theta} \neq \mathbf{0}$, so we would choose H_1 over H_0 if v_1 is sufficiently large. When Ω is positive semidefinite, the construction of v_1 above can be generalized by using the Moore–Penrose inverse of Ω . In this case, if $\text{rank}(\Omega) = r$, and we write $\Omega = X_1\Lambda_1X_1'$ and $\Omega^+ = X_1\Lambda_1^{-1}X_1'$, where the $m \times r$ matrix X_1 and the $r \times r$ diagonal matrix Λ_1 are defined as in the proof of Theorem 5.7, then $\mathbf{w} = \Lambda_1^{-1/2}X_1'\mathbf{t} \sim N_r(\Lambda_1^{-1/2}X_1'\boldsymbol{\theta}, I_r)$, because

$$\begin{aligned} \text{var}(\mathbf{w}) &= \Lambda_1^{-1/2}X_1'\{\text{var}(\mathbf{t})\}X_1\Lambda_1^{-1/2} \\ &= \Lambda_1^{-1/2}X_1'(X_1\Lambda_1X_1')X_1\Lambda_1^{-1/2} \\ &= I_r. \end{aligned}$$

Thus, because the w_i 's are independently distributed normal random variables,

$$v_2 = \mathbf{t}'\Omega^+\mathbf{t} = \mathbf{w}'\mathbf{w} = \sum_{i=1}^r w_i^2$$

has a chi-squared distribution, which is central if $\Lambda_1^{-1/2}X_1'\boldsymbol{\theta} = \mathbf{0}$, with r degrees of freedom.

5.4 THE MOORE–PENROSE INVERSE OF A MATRIX PRODUCT

If A and B each is an $m \times m$ nonsingular matrix, then it follows that $(AB)^{-1} = B^{-1}A^{-1}$. This property of the matrix inverse does not immediately generalize to the Moore–Penrose inverse of a matrix; that is, if A is $m \times p$ and B is $p \times n$, then we cannot, in general, be assured that $(AB)^+ = B^+A^+$. In this section, we look at some results regarding this sort of factorization of the Moore–Penrose inverse of a product.

Example 5.5 Here we look at a very simple example, given by Greville (1966), that illustrates a situation in which the factorization does not hold. Define the 2×1 vectors

$$\mathbf{a} = \begin{bmatrix} 1 \\ 0 \end{bmatrix}, \quad \mathbf{b} = \begin{bmatrix} 1 \\ 1 \end{bmatrix},$$

so that

$$\mathbf{a}^+ = (\mathbf{a}'\mathbf{a})^{-1}\mathbf{a}' = [1 \quad 0], \quad \mathbf{b}^+ = (\mathbf{b}'\mathbf{b})^{-1}\mathbf{b}' = [0.5 \quad 0.5].$$

Thus, we have

$$(a'b)^+ = (1)^+ = 1 \neq b^+a^{+'} = 0.5.$$

Actually, in the previous section, we have already given a few situations in which the identity $(AB)^+ = B^+A^+$ does hold. For example, in Theorem 5.3, we saw that

$$(A'A)^+ = A^+A^{+'} = A^+A^{+'}$$

and

$$(AA^+)^+ = AA^+ = (A^+)^+A^+.$$

Theorem 5.8 gives yet another situation.

Theorem 5.8 Let A be an $m \times n$ matrix, whereas P and Q are $h \times m$ and $n \times p$ matrices satisfying $P'P = I_m$ and $QQ' = I_n$. Then

$$(PAQ)^+ = Q^+A^+P^+ = Q^+A^+P^+.$$

The proof of Theorem 5.8, which we leave to the reader, simply involves the verification of conditions (5.1)–(5.4). Note that Theorem 5.7, regarding the Moore–Penrose inverse of a symmetric matrix, is a special case of Theorem 5.8.

Our next result gives a sufficient condition on the matrices A and B to guarantee that $(AB)^+ = B^+A^+$.

Theorem 5.9 Let A be an $m \times p$ matrix and B be a $p \times n$ matrix. If $\text{rank}(A) = \text{rank}(B) = p$, then $(AB)^+ = B^+A^+$.

Proof. Since A is full column rank and B is full row rank, we know from Theorem 5.3 that $A^+ = (A'A)^{-1}A'$ and $B^+ = B'(BB')^{-1}$. Consequently, we find that

$$\begin{aligned} ABB^+A^+AB &= ABB'(BB')^{-1}(A'A)^{-1}A'AB = AB, \\ B^+A^+ABB^+A^+ &= B'(BB')^{-1}(A'A)^{-1}A'ABB'(BB')^{-1}(A'A)^{-1}A' \\ &= B'(BB')^{-1}(A'A)^{-1}A' = B^+A^+, \end{aligned}$$

so conditions (5.1) and (5.2) are satisfied. In addition,

$$\begin{aligned} ABB^+A^+ &= ABB'(BB')^{-1}(A'A)^{-1}A' \\ &= A(A'A)^{-1}A', \\ B^+A^+AB &= B'(BB')^{-1}(A'A)^{-1}A'AB \\ &= B'(BB')^{-1}B \end{aligned}$$

are symmetric, so B^+A^+ is the Moore–Penrose inverse of AB . □

Although Theorem 5.9 is useful, its major drawback is that it only gives a sufficient condition for the factorization of $(AB)^+$. The following result, due to Greville (1966), gives several necessary and sufficient conditions for this factorization to hold.

Theorem 5.10 Let A be an $m \times p$ matrix and B be a $p \times n$ matrix. Then each of the following conditions are necessary and sufficient for $(AB)^+ = B^+A^+$:

- (a) $A^+ABB'A' = BB'A'$ and $BB^+A'AB = A'AB$.
- (b) A^+ABB' and $A'ABB^+$ are symmetric matrices.
- (c) $A^+ABB'A'ABB^+ = BB'A'A$.
- (d) $A^+AB = B(AB)^+AB$ and $BB^+A' = A'AB(AB)^+$.

Proof. We will prove that the conditions given in (a) are necessary and sufficient; the proofs for (b)–(d) will be left to the reader as an exercise. First assume that the conditions of (a) hold. Premultiplying the first identity by B^+ while postmultiplying by $(AB)^{+'}$ yields

$$B^+A^+AB(AB)'(AB)^{+'} = B^+BB'A'(AB)^{+'}. \quad (5.8)$$

Now for any matrix C ,

$$\begin{aligned} C^+CC' &= (C^+C)'C' = C'C^{+'}C' \\ &= C'C'^+C' = C'. \end{aligned} \quad (5.9)$$

Using this identity, when $C = B$, on the right-hand side of (5.8) and its transpose on the left-hand side, when $C = AB$, we obtain the equation

$$B^+A^+AB = (AB)'(AB)^{+'},$$

which, because of condition (5.4), is equivalent to

$$B^+A^+AB = (AB)^+(AB) = P_{R((AB)^+)}. \quad (5.10)$$

The final equality in (5.10) follows from the definition of the Moore–Penrose inverse in term of projection matrices, as given in Definition 5.2. In a similar fashion, if we take the transpose of the second identity in (a), which yields

$$B'A'ABB^+ = B'A'A$$

and premultiply this by $(AB)^{+'}$ and postmultiply this by A^+ , then, after simplifying by using (5.9) on the left-hand side with $C = (AB)'$ and the transpose of (5.9) on the right-hand side with $C = A'$, we obtain the equation

$$ABB^+A^+ = (AB)(AB)^+ = P_{R(AB)}. \quad (5.11)$$

However, from Definition 5.2, $(AB)^+$ is the only matrix satisfying both (5.10) and (5.11). Consequently, we must have $(AB)^+ = B^+A^+$. Conversely, now suppose that $(AB)^+ = B^+A^+$. Applying this equation in (5.9), when $C = AB$, gives

$$(AB)' = B^+A^+(AB)(AB)'.$$

Premultiplying this by $ABB'B$, we obtain

$$ABB'BB'A' = ABB'BB^+A^+ABB'A',$$

which, after using the transpose of (5.9) with $C = B'$ and then rearranging, simplifies to

$$ABB'(I_p - A^+A)BB'A' = (0).$$

Note that because $D = (I_p - A^+A)$ is symmetric and idempotent, the equation above is in the form $E'D'DE = (0)$, where $E = BB'A'$. This then implies that $ED = (0)$; that is,

$$(I_p - A^+A)BB'A' = (0),$$

which is equivalent to the first identity in (a). In a similar fashion, using $(AB)^+ = B^+A^+$ in (5.9) with $C = (AB)'$ yields

$$AB = A^{+'}B^{+'}B'A'AB.$$

This, when premultiplied by $B'A'AA'$, can be simplified to an equation that is equivalent to the second identity of (a). \square

Our next step is to find a general expression for $(AB)^+$ that holds for all A and B for which the product AB is defined. Our approach involves transforming A to a matrix A_1 and transforming B to B_1 , such that $AB = A_1B_1$ and $(A_1B_1)^+ = B_1^+A_1^+$. The result, due to Cline (1964a), is given in Theorem 5.11.

Theorem 5.11 Let A be an $m \times p$ matrix and B be a $p \times n$ matrix. If we define $B_1 = A^+AB$ and $A_1 = AB_1B_1^+$, then $AB = A_1B_1$ and $(AB)^+ = B_1^+A_1^+$.

Proof. Note that

$$AB = AA^+AB = AB_1 = AB_1B_1^+B_1 = A_1B_1,$$

so the first result holds. To verify the second statement, we will show that the two conditions given in Theorem 5.10(a) are satisfied for A_1 and B_1 . First note that

$$\begin{aligned} A^+A_1 &= A^+AB_1B_1^+ = A^+A(A^+AB)B_1^+ \\ &= A^+ABB_1^+ = B_1B_1^+ \end{aligned} \tag{5.12}$$

and

$$\begin{aligned} A_1^+ A_1 &= A_1^+ A B_1 B_1^+ = A_1^+ A (B_1 B_1^+ B_1) B_1^+ \\ &= A_1^+ A_1 B_1 B_1^+. \end{aligned} \quad (5.13)$$

Taking the transpose of (5.13) and using (5.12), along with conditions (5.3) and (5.4), we get

$$A_1^+ A_1 = B_1 B_1^+ A_1^+ A_1 = A^+ A_1 A_1^+ A_1 = A^+ A_1 = B_1 B_1^+,$$

and so

$$A_1^+ A_1 B_1 B_1^+ A_1' = B_1 B_1^+ B_1 B_1^+ A_1' = B_1 B_1^+ A_1',$$

which is the first identity in Theorem 5.10(a). The second identity can be obtained by noting that

$$\begin{aligned} A_1' &= (A B_1 B_1^+)' = (A B_1 B_1^+ B_1 B_1^+)' \\ &= (A_1 B_1 B_1^+)' = B_1 B_1^+ A_1', \end{aligned}$$

and then postmultiplying this identity by $A_1 B_1$. □

Note that in Theorem 5.11, B was transformed to B_1 by the projection matrix of the range space of A^+ , whereas A was transformed to A_1 by the projection matrix of the range space of B_1 and not that of B . Our next result indicates that the range space of B can be used instead of that of B_1 , if we do not insist that $AB = A_1 B_1$. A proof of this result can be found in Campbell and Meyer (1979).

Theorem 5.12 Let A be an $m \times p$ matrix and B be a $p \times n$ matrix. If we define $B_1 = A^+ A B$ and $A_1 = A B B^+$, then $(AB)^+ = B_1^+ A_1^+$.

5.5 THE MOORE–PENROSE INVERSE OF PARTITIONED MATRICES

Suppose that the $m \times n$ matrix A has been partitioned as $A = [U \ V]$, where U is $m \times n_1$ and V is $m \times n_2$. In some situations, it may be useful to have an expression for A^+ in terms of the submatrices, U and V . We begin with the general case, in which no assumptions can be made regarding U and V .

Theorem 5.13 Let the $m \times n$ matrix A be partitioned as $A = [U \ V]$, where U is $m \times n_1$, V is $m \times n_2$, and $n = n_1 + n_2$. Then

$$A^+ = \begin{bmatrix} U^+ - U^+ V (C^+ + W) \\ C^+ + W \end{bmatrix},$$

where $C = (I_m - U U^+) V$, $M = \{I_{n_2} + (I_{n_2} - C^+ C) V' U^+ U^+ V (I_{n_2} - C^+ C)\}^{-1}$, and $W = (I_{n_2} - C^+ C) M V' U^+ U^+ (I_m - V C^+)$.

Proof. Partition A^+ as

$$A^+ = \begin{bmatrix} X \\ Y \end{bmatrix},$$

so that

$$AA^+ = UX + VY \quad (5.14)$$

and

$$A^+A = \begin{bmatrix} XU & XV \\ YU & YV \end{bmatrix}. \quad (5.15)$$

Since $AA^+A = A$, we have

$$AA^+U = U, \quad (5.16)$$

$$AA^+V = V. \quad (5.17)$$

Transposing (5.16) and then premultiplying by $(U'U)^+$, we get

$$U^+AA^+ = U^+, \quad (5.18)$$

because $(U'U)^+U' = U^+$ by Theorem 5.3(g). Also from Theorem 5.3(e) and (g), we have $A^+ = A'A^+A$, so that

$$X = U'X'X + U'Y'Y,$$

which leads to

$$\begin{aligned} U^+UX &= U^+UU'X'X + U^+UU'Y'Y \\ &= U'X'X + U'Y'Y \\ &= X. \end{aligned}$$

Thus, premultiplying (5.14) by U^+ and using (5.18), we find that $U^+ = X + U^+VY$, which implies that

$$A^+ = \begin{bmatrix} U^+ - U^+VY \\ Y \end{bmatrix}. \quad (5.19)$$

Consequently, the proof will be complete if we can show that $Y = C^+ + W$. Since $U^+C = (0)$, it follows (Problem 5.14) that

$$C^+U = (0), \quad (5.20)$$

and using (5.16) and (5.17), we get $AA^+C = C$, or equivalently $C' = C'AA^+$. This last identity, when premultiplied by $(C'C)^+$, yields $C^+ = C^+AA^+$. Thus, using (5.19), we have

$$AA^+ = UU^+ + (I_m - UU^+)VY = UU^+ + CY,$$

and when this identity is premultiplied by C^+ , it reduces to $C^+ = C^+CY$, so that

$$CY = CC^+. \quad (5.21)$$

Also, $C = CC^+C = CC^+(V - UU^+V) = CC^+V$, which implies

$$C^+V = C^+C. \quad (5.22)$$

Thus, $CYV = CC^+V = CC^+C = C$, and from (5.15) we know that YV is symmetric, so $YVC' = C'$ or when postmultiplying this last identity by $(CC')^+$,

$$YVC^+ = C^+. \quad (5.23)$$

Using the expression for A^+ given in (5.19) and the identity $A^+AA^+ = A^+$, we find that

$$YUU^+ + YCY = Y$$

or

$$YUU^+ + YCC^+ = Y \quad (5.24)$$

because of (5.21). The symmetry condition $(A^+A)' = A^+A$ yields the equations

$$U^+VYU = (U^+VYU)', \quad (5.25)$$

$$(YU)' = U^+V(I_{n_2} - YV). \quad (5.26)$$

Now from (5.26) and the definition of C ,

$$\begin{aligned} (YU)' &= U^+V\{I_{n_2} - Y(UU^+V + C)\} \\ &= U^+V - U^+VYUU^+V - U^+VYC, \end{aligned}$$

and because by (5.20) and (5.21)

$$(I_{n_2} - C^+C)YU = YU - C^+CYU = YU - C^+CC^+U = YU,$$

it follows that

$$\begin{aligned} (YU)' &= \{(I_{n_2} - C^+C)YU\}' = (YU)'(I_{n_2} - C^+C) \\ &= (U^+V - U^+VYUU^+V - U^+VYC)(I_{n_2} - C^+C) \\ &= (U^+V - U^+VYUU^+V)(I_{n_2} - C^+C). \end{aligned}$$

Transposing this last equation and using (5.25), we get

$$\begin{aligned} YU &= (I_{n_2} - C^+C)V'U^{+'} - (I_{n_2} - C^+C)V'U^{+'}U^+VYU \\ &= (I_{n_2} - C^+C)V'U^{+'} - (I_{n_2} - C^+C)V'U^{+'}U^+V(I_{n_2} - C^+C)YU, \end{aligned}$$

which leads to

$$BYU = (I_{n_2} - C^+C)V'U^{+'}, \quad (5.27)$$

where $B = I_{n_2} + (I_{n_2} - C^+C)V'U^{+'}U^+V(I_{n_2} - C^+C)$. Postmultiplying (5.27) by $U^+(I_m - VC^+)$ and then using (5.22), (5.23) and (5.24), we have

$$B(Y - C^+) = (I_{n_2} - C^+C)V'U^{+'}U^+(I_m - VC^+).$$

Since B is the sum of I_{n_2} and a nonnegative definite matrix, it must be positive definite and, hence, nonsingular. Thus, our previous equation can be re-expressed as

$$\begin{aligned} Y &= C^+ + B^{-1}(I_{n_2} - C^+C)V'U^{+'}U^+(I_m - VC^+) \\ &= C^+ + (I_{n_2} - C^+C)B^{-1}V'U^{+'}U^+(I_m - VC^+) \\ &= C^+ + W, \end{aligned}$$

where we have used the fact that B^{-1} and $(I_{n_2} - C^+C)$ commute because B and $(I_{n_2} - C^+C)$ commute, and $B^{-1} = M$. This completes the proof. \square

The proofs of the following consequences of Theorem 5.13 can be found in Cline (1964b), Boullion and Odell (1971), or Pringle and Rayner (1971).

Corollary 5.13.1 Let A and C be defined as in Theorem 5.13, and let $K = (I_{n_2} + V'U^{+'}U^+V)^{-1}$. Then

(a) $C^+CV'U^{+'}U^+V = (0)$ if and only if

$$A^+ = \begin{bmatrix} U^+ - U^+VKV'U^{+'}U^+ \\ C^+ + KV'U^{+'}U^+ \end{bmatrix},$$

(b) $C = (0)$ if and only if

$$A^+ = \begin{bmatrix} U^+ - U^+VKV'U^{+'}U^+ \\ KV'U^{+'}U^+ \end{bmatrix},$$

(c) $C^+CV'U^+U^+V = V'U^+U^+V$ if and only if

$$A^+ = \begin{bmatrix} U^+ - U^+VC^+ \\ C^+ \end{bmatrix},$$

(d) $U'V = (0)$ if and only if

$$A^+ = \begin{bmatrix} U^+ \\ V^+ \end{bmatrix}.$$

Our final theorem involves the Moore–Penrose inverse of a partitioned matrix that has the block diagonal form. This result can be easily proven by simply verifying that the conditions of the Moore–Penrose inverse are satisfied.

Theorem 5.14 Let the $m \times n$ matrix A be given by

$$A = \begin{bmatrix} A_{11} & (0) & \cdots & (0) \\ (0) & A_{22} & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & A_{rr} \end{bmatrix},$$

where A_{ii} is $m_i \times n_i$, $m_1 + \cdots + m_r = m$, and $n_1 + \cdots + n_r = n$. Then

$$A^+ = \begin{bmatrix} A_{11}^+ & (0) & \cdots & (0) \\ (0) & A_{22}^+ & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & A_{rr}^+ \end{bmatrix}.$$

Some additional results for the generalized inverse of a matrix A that is partitioned into a 2×2 form will be given in Chapter 7.

5.6 THE MOORE–PENROSE INVERSE OF A SUM

Theorem 1.9 gave an expression for $(A + CBD)^{-1}$, when the matrices A , B , and $A + CBD$ are all square and nonsingular. Although a generalization of this formula to the case of a Moore–Penrose inverse is not available, there are some specialized results for the Moore–Penrose inverse of a sum of matrices. Some of these results are presented in this section. The proofs of our first two results use the results of the previous section regarding partitioned matrices. These proofs can be found in Cline (1965) or Boullion and Odell (1971).

Theorem 5.15 Let U be an $m \times n_1$ matrix and V be an $m \times n_2$ matrix. Then

$$(UU' + VV')^+ = (I_m - C^+V')U^+KU^+(I_m - VC^+) + (CC')^+,$$

where $K = I_{n_1} - U^+V(I_{n_2} - C^+C)M(U^+V)'$ and C and M are defined as in Theorem 5.13.

Theorem 5.16 Suppose U and V are both $m \times n$ matrices. If $UV' = (0)$, then

$$(U + V)^+ = U^+ + (I_n - U^+V)(C^+ + W),$$

where C and W are as given in Theorem 5.13.

Theorem 5.16 gives an expression for $(U + V)^+$ that holds when the rows of U are orthogonal to the rows of V . If, in addition, the columns of U are orthogonal to the columns of V , this expression greatly simplifies. This special case is summarized in Theorem 5.17.

Theorem 5.17 If U and V are $m \times n$ matrices satisfying $UV' = (0)$ and $U'V = (0)$, then

$$(U + V)^+ = U^+ + V^+.$$

Proof. Using Theorem 5.3(g), we find that

$$U^+V = (U'U)^+U'V = (0)$$

and

$$VV^+ = VU'(UU')^+ = \{(UU')^+U'V'\}' = (0).$$

Similarly, we have $V^+U = (0)$ and $UV^+ = (0)$. As a result,

$$(U + V)(U^+ + V^+) = UU^+ + VV^+, \quad (5.28)$$

$$(U^+ + V^+)(U + V) = U^+U + V^+V, \quad (5.29)$$

which are both symmetric, so that conditions (5.3) and (5.4) are satisfied. Postmultiplying (5.28) by $(U + V)$ and (5.29) by $(U^+ + V^+)$ yields conditions (5.1) and (5.2), so the result follows. \square

Theorem 5.17 can be easily generalized to more than two matrices.

Corollary 5.17.1 Let U_1, \dots, U_k be $m \times n$ matrices satisfying $U_iU_j' = (0)$ and $U_i'U_j = (0)$ for all $i \neq j$. Then

$$(U_1 + \dots + U_k)^+ = U_1^+ + \dots + U_k^+.$$

We saw in Corollary 1.9.2 that if A and $A + \mathbf{cd}'$ are nonsingular matrices, then

$$(A + \mathbf{cd}')^{-1} = A^{-1} - \frac{A^{-1}\mathbf{cd}'A^{-1}}{1 + \mathbf{d}'A^{-1}\mathbf{c}}.$$

Our final theorem gives a generalization of this result to the case in which $A + \mathbf{cd}'$ is singular and A is symmetric.

Theorem 5.18 Let A be an $m \times m$ nonsingular symmetric matrix and \mathbf{c} and \mathbf{d} be $m \times 1$ vectors. Then $A + \mathbf{cd}'$ is singular if and only if $1 + \mathbf{d}'A^{-1}\mathbf{c} = 0$, and if $A + \mathbf{cd}'$ is singular,

$$(A + \mathbf{cd}')^+ = (I_m - \mathbf{yy}^+)A^{-1}(I_m - \mathbf{xx}^+),$$

where $\mathbf{x} = A^{-1}\mathbf{d}$ and $\mathbf{y} = A^{-1}\mathbf{c}$.

Proof. Our proof follows that of Trenkler (2000). From Theorem 3.6, $|A + \mathbf{cd}'| = |A|(1 + \mathbf{d}'A^{-1}\mathbf{c})$, and so the stated necessary and sufficient condition for the singularity of $A + \mathbf{cd}'$ follows. Now if $1 + \mathbf{d}'A^{-1}\mathbf{c} = 0$, it follows that $(A + \mathbf{cd}')\mathbf{y} = \mathbf{c} + \mathbf{c}(\mathbf{d}'A^{-1}\mathbf{c}) = \mathbf{c} - \mathbf{c} = \mathbf{0}$, which implies $(A + \mathbf{cd}')\mathbf{yy}^+ = (\mathbf{0})$, or equivalently

$$(A + \mathbf{cd}')(I_m - \mathbf{yy}^+) = A + \mathbf{cd}'. \quad (5.30)$$

In a similar fashion, because A is symmetric we can show that

$$(I_m - \mathbf{xx}^+)(A + \mathbf{cd}') = A + \mathbf{cd}'. \quad (5.31)$$

Thus, using (5.31), we get

$$\begin{aligned} (I_m - \mathbf{yy}^+)A^{-1}(I_m - \mathbf{xx}^+)(A + \mathbf{cd}') &= (I_m - \mathbf{yy}^+)A^{-1}(A + \mathbf{cd}') \\ &= (I_m - \mathbf{yy}^+)(I_m + \mathbf{yd}') \\ &= (I_m - \mathbf{yy}^+), \end{aligned} \quad (5.32)$$

whereas an application of (5.30) confirms that

$$\begin{aligned} (A + \mathbf{cd}')(I_m - \mathbf{yy}^+)A^{-1}(I_m - \mathbf{xx}^+) &= (A + \mathbf{cd}')A^{-1}(I_m - \mathbf{xx}^+) \\ &= (I_m + \mathbf{cx}')(I_m - \mathbf{xx}^+) \\ &= (I_m - \mathbf{xx}^+). \end{aligned} \quad (5.33)$$

This establishes conditions (5.3) and (5.4) of a Moore–Penrose inverse. Condition (5.1) follows by premultiplying (5.32) by $(A + \mathbf{cd}')$ and then applying (5.30), whereas condition (5.2) is obtained by postmultiplying (5.33) by $(A + \mathbf{cd}')$ and then applying (5.31). This completes the proof. \square

5.7 THE CONTINUITY OF THE MOORE–PENROSE INVERSE

It is very useful to establish the continuity of a function because continuous functions enjoy many nice properties. In this section, we will give conditions under which the elements of A^+ are continuous functions of the elements of A . However, before doing so, let us first consider the determinant of a square matrix A and the inverse of a nonsingular matrix A . Recall that the determinant of an $m \times m$ matrix A can be expressed as the sum of terms, where each term is $+1$ or -1 times the product of m of the elements of A . Thus, because of the continuity of sums and the continuity of scalar products, we immediately have the following.

Theorem 5.19 Let A be an $m \times m$ matrix. Then the determinant of A , $|A|$, is a continuous function of the elements of A .

Suppose that A is an $m \times m$ nonsingular matrix so that $|A| \neq 0$. Recall that the inverse of A can be expressed as

$$A^{-1} = |A|^{-1} A_{\#}, \quad (5.34)$$

where $A_{\#}$ is the adjoint matrix of A . If A_1, A_2, \dots is a sequence of matrices such that $A_i \rightarrow A$ as $i \rightarrow \infty$, then, because of the continuity of the determinant function, $|A_i| \rightarrow |A|$, and so an N must exist, such that $|A_i| \neq 0$ for all $i > N$. Since each element of an adjoint matrix is $+1$ or -1 times a determinant, it also follows from the continuity of the determinant function that if $A_{i\#}$ is the adjoint matrix of A_i , then $A_{i\#} \rightarrow A_{\#}$ as $i \rightarrow \infty$. As a result, (5.34) has allowed us to establish the following.

Theorem 5.20 Let A be an $m \times m$ nonsingular matrix. Then the inverse of A , A^{-1} , is a continuous function of the elements of A .

An alternative way of establishing the continuity of A^{-1} is to show directly that if $A_i \rightarrow A$, then for some matrix norm, $\|A^{-1} - A_i^{-1}\| \rightarrow 0$. Without loss of generality, we assume our norm satisfies the identity $\|I_m\| = 1$. Let $B_i = A_i - A$ so that $A_i = A + B_i = A(I_m + A^{-1}B_i)$. Now since $B_i \rightarrow (0)$, there exists an integer N such that $\|B_i\| < 1/\|A^{-1}\|$ for all $i > N$. Thus, following the derivation given in Example 4.14, we find that for $i > N$,

$$\|A^{-1} - (A + B_i)^{-1}\| \leq \frac{\|A^{-1}\|^2 \|B_i\|}{1 - \|A^{-1}\| \|B_i\|}.$$

From this it immediately follows that $\|A^{-1} - A_i^{-1}\| = \|A^{-1} - (A + B_i)^{-1}\| \rightarrow 0$ since $\|B_i\| \rightarrow 0$.

The continuity of the Moore–Penrose inverse is not as straightforward as the continuity of the inverse of a nonsingular matrix. If A is an $m \times n$ matrix and A_1, A_2, \dots is an arbitrary sequence of $m \times n$ matrices satisfying $A_i \rightarrow A$ as $i \rightarrow \infty$, then we are not assured that $A_i^+ \rightarrow A^+$. A simple example will illustrate the potential problem.

Example 5.6 Consider the sequence of 2×2 matrices A_1, A_2, \dots , where

$$A_i = \begin{bmatrix} 1/i & 0 \\ 0 & 1 \end{bmatrix}.$$

Clearly, $A_i \rightarrow A$, where

$$A = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

However, note that $\text{rank}(A) = 1$, whereas $\text{rank}(A_i) = 2$ for all i . For this reason, we do not have $A_i^+ \rightarrow A^+$. In fact,

$$A_i^+ = \begin{bmatrix} i & 0 \\ 0 & 1 \end{bmatrix}$$

does not converge to anything because its $(1, 1)$ th element, i , goes to ∞ . On the other hand,

$$A^+ = \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix}.$$

If we have a sequence of matrices A_1, A_2, \dots for which $\text{rank}(A_i) = \text{rank}(A)$ for all i larger than some integer, say N , then we will not encounter the difficulty observed in Example 5.6; that is, as A_i gets closer to A , A_i^+ will get closer to A^+ . This continuity property of A^+ is summarized in Theorem 5.21. A proof of this important result can be found in Penrose (1955) or Campbell and Meyer (1979).

Theorem 5.21 Let A be an $m \times n$ matrix and A_1, A_2, \dots be a sequence of $m \times n$ matrices, such that $A_i \rightarrow A$ as $i \rightarrow \infty$. Then

$$A_i^+ \rightarrow A^+ \quad \text{as } i \rightarrow \infty$$

if and only if an integer N exists, such that

$$\text{rank}(A_i) = \text{rank}(A) \quad \text{for all } i > N.$$

Example 5.7 The conditions for the continuity of the Moore–Penrose inverse have important implications in estimation and hypothesis testing problems. In particular, in this example, we will discuss a property, referred to as consistency, that some estimators possess. An estimator t , computed from a sample of size n , is said to be a consistent estimator of a parameter θ if t converges in probability to θ , that is, if

$$\lim_{n \rightarrow \infty} P(|t - \theta| \geq \epsilon) = 0,$$

for any $\epsilon > 0$. An important result associated with the property of consistency is that continuous functions of consistent estimators are consistent; that is, if t is a consistent estimator of θ , and $g(t)$ is a continuous function of t , then $g(t)$ is a consistent estimator of $g(\theta)$. We will now apply some of these ideas to a situation involving the estimation of the Moore–Penrose inverse of a matrix of parameters. For instance, let Ω be an $m \times m$ positive semidefinite covariance matrix having rank $r < m$. Suppose that the elements of the matrix Ω are unknown and are, therefore, to be estimated. Suppose, in addition, that our sample estimate of Ω , which we will denote by $\hat{\Omega}$, is positive definite with probability one, so that $\text{rank}(\hat{\Omega}) = m$ with probability one, and $\hat{\Omega}$ is a consistent estimator of Ω ; that is, each element of $\hat{\Omega}$ is a consistent estimator of the corresponding element of Ω . However, because $\text{rank}(\Omega) = r < m$, $\hat{\Omega}^+$ is not a consistent estimator Ω^+ . Intuitively, the problem here is obvious. If $\hat{\Omega} = X\Lambda X'$ is the spectral decomposition of $\hat{\Omega}$ so that $\hat{\Omega}^+ = \hat{\Omega}^{-1} = X\Lambda^{-1}X'$, then the consistency of $\hat{\Omega}$ is implying that as n increases, the $m - r$ smallest diagonal elements of Λ are converging to zero, whereas the $m - r$ largest diagonal elements of Λ^{-1} are increasing without bound. The difficulty here can be easily avoided if the value of r is known. In this case, $\hat{\Omega}$ can be adjusted to yield an estimator of Ω having rank r . For example, if $\hat{\Omega}$ has eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$ and corresponding normalized eigenvectors $\mathbf{x}_1, \dots, \mathbf{x}_m$, and P_r is the eigenprojection

$$P_r = \sum_{i=1}^r \mathbf{x}_i \mathbf{x}_i',$$

then

$$\hat{\Omega}_* = P_r \hat{\Omega} P_r = \sum_{i=1}^r \lambda_i \mathbf{x}_i \mathbf{x}_i'$$

will be an estimator of Ω having rank of r . It can be shown then that, because of the continuity of eigenprojections, $\hat{\Omega}_*$ is also a consistent estimator of Ω . More importantly, because $\text{rank}(\hat{\Omega}_*) = \text{rank}(\Omega) = r$, Theorem 5.21 guarantees that $\hat{\Omega}_*^+$ is a consistent estimator of Ω^+ .

5.8 SOME OTHER GENERALIZED INVERSES

The Moore–Penrose inverse is just one of many generalized inverses that have been developed in recent years. In this section, we will briefly discuss two other generalized inverses that have applications in statistics. Both of these inverses can be defined by applying some of the four conditions, (5.1)–(5.4), or, for simplicity, 1–4, of the Moore–Penrose inverse. In fact, we can define a different class of inverses corresponding to each different subset of the conditions 1–4 that the inverse must satisfy.

Definition 5.3 For any $m \times n$ matrix A , let the $n \times m$ matrix denoted $A^{(i_1, \dots, i_r)}$ be any matrix satisfying conditions i_1, \dots, i_r from among the four conditions 1–4; $A^{(i_1, \dots, i_r)}$ will be called a $\{i_1, \dots, i_r\}$ -inverse of A .

Thus, the Moore–Penrose inverse of A is the $\{1, 2, 3, 4\}$ -inverse of A ; that is, $A^+ = A^{(1,2,3,4)}$. Note that for any proper subset $\{i_1, \dots, i_r\}$ of $\{1, 2, 3, 4\}$, A^+ will also be a $\{i_1, \dots, i_r\}$ -inverse of A , but it may not be the only one. Since in many cases, there are many different $\{i_1, \dots, i_r\}$ -inverses of A , it may be easier to compute a $\{i_1, \dots, i_r\}$ -inverse of A than to compute the Moore–Penrose inverse. The rest of this section will be devoted to the $\{1\}$ -inverse of A and the $\{1, 3\}$ -inverse of A , which have special applications that will be discussed in Chapter 6. Discussion of other useful $\{i_1, \dots, i_r\}$ -inverses can be found in Ben-Israel and Greville (2003), Campbell and Meyer (1979), and Rao and Mitra (1971).

In Chapter 6, we will see that in solving systems of linear equations, we will only need an inverse matrix satisfying the first condition of the four Moore–Penrose conditions. We will refer to any such $\{1\}$ -inverse of A as simply a generalized inverse of A , and we will write it using the fairly common notation, A^- ; that is, $A^{(1)} = A^-$. One useful way of expressing a generalized inverse of a matrix A applies the singular value decomposition of A . The following result, which is stated for a matrix A having less than full rank, can easily be modified for matrices having full row rank or full column rank.

Theorem 5.22 Suppose that the $m \times n$ matrix A has rank $r > 0$ and the singular value decomposition given by

$$A = P \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q',$$

where P and Q are $m \times m$ and $n \times n$ orthogonal matrices, respectively, and Δ is an $r \times r$ nonsingular diagonal matrix. Let

$$B = Q \begin{bmatrix} \Delta^{-1} & E \\ F & G \end{bmatrix} P',$$

where E is $r \times (m - r)$, F is $(n - r) \times r$, and G is $(n - r) \times (m - r)$. Then for all choices of E , F , and G , B is a generalized inverse of A , and any generalized inverse of A can be expressed in the form of B for some E , F , and G .

Proof. Note that

$$\begin{aligned} ABA &= P \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q' Q \begin{bmatrix} \Delta^{-1} & E \\ F & G \end{bmatrix} P' P \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q' \\ &= P \begin{bmatrix} \Delta \Delta^{-1} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q' \\ &= P \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q' = A, \end{aligned}$$

and so the matrix B is a generalized inverse of A regardless of the choice of E , F , and G . On the other hand, if we write $Q = [Q_1 \ Q_2]$, $P = [P_1 \ P_2]$, where

Q_1 is $n \times r$ and P_1 is $m \times r$, then, because $PP' = I_m$, $QQ' = I_n$, any generalized inverse B , of A , can be expressed as

$$\begin{aligned} B &= QQ'BP P' = Q \begin{bmatrix} Q'_1 \\ Q'_2 \end{bmatrix} B \begin{bmatrix} P_1 & P_2 \end{bmatrix} P' \\ &= Q \begin{bmatrix} Q'_1 B P_1 & Q'_1 B P_2 \\ Q'_2 B P_1 & Q'_2 B P_2 \end{bmatrix} P', \end{aligned}$$

which is in the required form if we can show that $Q'_1 B P_1 = \Delta^{-1}$. Since B is a generalized inverse of A , $ABA = A$, or equivalently,

$$(P' A Q)(Q' B P)(P' A Q) = P' A Q.$$

Writing this last identity in partitioned form and equating the $(1, 1)$ th submatrices on both sides, we find that

$$\Delta Q'_1 B P_1 \Delta = \Delta,$$

from which it immediately follows that $Q'_1 B P_1 = \Delta^{-1}$, and so the proof is complete. \square

When A is an $m \times m$ nonsingular matrix, the matrix B in Theorem 5.22 simplifies to $B = Q\Delta^{-1}P'$, where Δ , P , and Q are now all $m \times m$ matrices. In other words, an immediate consequence of Theorem 5.22 is that A^{-1} is the only generalized inverse of A when A is square and nonsingular.

Example 5.8 The 4×3 matrix

$$A = \begin{bmatrix} 1 & 0 & 0.5 \\ 1 & 0 & 0.5 \\ 0 & -1 & -0.5 \\ 0 & -1 & -0.5 \end{bmatrix}$$

has rank $r = 2$ and singular value decomposition with

$$P = \frac{1}{2} \begin{bmatrix} 1 & 1 & 1 & -1 \\ 1 & 1 & -1 & 1 \\ 1 & -1 & 1 & 1 \\ 1 & -1 & -1 & -1 \end{bmatrix}, \quad Q' = \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{6} & 1/\sqrt{6} & -2/\sqrt{6} \end{bmatrix},$$

and

$$\Delta = \begin{bmatrix} \sqrt{2} & 0 \\ 0 & \sqrt{3} \end{bmatrix}.$$

If we take E , F , and G as null matrices and use the equation for B given in Theorem 5.22, we obtain as a generalized inverse of A the matrix

$$\frac{1}{12} \begin{bmatrix} 5 & 5 & 1 & 1 \\ -1 & -1 & -5 & -5 \\ 2 & 2 & -2 & -2 \end{bmatrix}.$$

Actually, from the proof of Theorem 5.1, we know that the matrix above is the Moore–Penrose inverse. Different generalized inverses of A may be constructed through different choices of E , F , and G ; for example, if we again take E and F as null matrices but now use

$$G = [1/\sqrt{6} \quad 0],$$

then we obtain the generalized inverse

$$\frac{1}{6} \begin{bmatrix} 3 & 2 & 1 & 0 \\ 0 & -1 & -2 & -3 \\ 0 & 2 & -2 & 0 \end{bmatrix}.$$

Note that this matrix has rank 3, whereas the Moore–Penrose inverse has its rank equal to that of A , which is 2.

Theorem 5.23 summarizes some of the basic properties of $\{1\}$ -inverses.

Theorem 5.23 Let A be an $m \times n$ matrix, and let A^- be a generalized inverse of A . Then

- (a) $A^{-'}$ is a generalized inverse of A' ,
- (b) if α is a nonzero scalar, $\alpha^{-1}A^-$ is a generalized inverse of αA ,
- (c) if A is square and nonsingular, $A^- = A^{-1}$ uniquely,
- (d) if B and C are nonsingular, $C^{-1}A^-B^{-1}$ is a generalized inverse of BAC ,
- (e) $\text{rank}(A) = \text{rank}(AA^-) = \text{rank}(A^-A) \leq \text{rank}(A^-)$,
- (f) $\text{rank}(A) = m$ if and only if $AA^- = I_m$,
- (g) $\text{rank}(A) = n$ if and only if $A^-A = I_n$,
- (h) if $m = n$ and A is symmetric, there exists a generalized inverse A^- that is symmetric,
- (i) if $m = n$ and A is nonnegative definite, there exists a generalized inverse A^- that is nonnegative definite.

Proof. Properties (a)–(d) are easily proven by simply verifying that the one condition of a generalized inverse holds. To prove (e), note that because $A = AA^-A$, we can use Theorem 2.8 to get

$$\text{rank}(A) = \text{rank}(AA^-A) \leq \text{rank}(AA^-) \leq \text{rank}(A)$$

and

$$\text{rank}(A) = \text{rank}(AA^-A) \leq \text{rank}(A^-A) \leq \text{rank}(A),$$

so that $\text{rank}(A) = \text{rank}(AA^-) = \text{rank}(A^-A)$. In addition,

$$\text{rank}(A) = \text{rank}(AA^-A) \leq \text{rank}(A^-A) \leq \text{rank}(A^-),$$

so the result follows. It follows from (e) that $\text{rank}(A) = m$ if and only if AA^- is nonsingular. Premultiplying the equation

$$(AA^-)^2 = (AA^-A)A^- = AA^-$$

by $(AA^-)^{-1}$ yields (f). Similarly, $\text{rank}(A) = n$ if and only if A^-A is nonsingular, and so premultiplying

$$(A^-A)^2 = A^-(AA^-A) = A^-A$$

by $(A^-A)^{-1}$ gives (g). Properties (h) and (i) follow immediately from the fact that A^+ is a generalized inverse of A . \square

Example 5.9 Some of the properties possessed by the Moore–Penrose inverse do not carry over to the $\{1\}$ -inverse. For instance, we have seen that A is the Moore–Penrose inverse of A^+ ; that is, $(A^+)^+ = A$. However, in general, we are not guaranteed that A is a generalized inverse of A^- , where A^- is an arbitrary generalized inverse of A . For example, consider the diagonal matrix $A = \text{diag}(0, 2, 4)$. One choice of a generalized inverse of A is $A^- = \text{diag}(1, 0.5, 0.25)$. Here A^- is nonsingular, so it has only one generalized inverse, namely, $(A^-)^{-1} = \text{diag}(1, 2, 4)$, and, thus, A is not a generalized inverse of $A^- = \text{diag}(1, 0.5, 0.25)$.

All generalized inverses of a matrix A can be expressed in terms of any one particular generalized inverse. This relationship is given below.

Theorem 5.24 Let A^- be any generalized inverse of the $m \times n$ matrix A . Then for any $n \times m$ matrix C ,

$$A^- + C - A^-ACAA^-$$

is a generalized inverse of A , and each generalized inverse of A can be expressed in this form for some C .

Proof. Since $AA^-A = A$,

$$\begin{aligned} A(A^- + C - A^-ACAA^-)A &= AA^-A + ACA - AA^-ACAA^-A \\ &= A + ACA - ACA = A, \end{aligned}$$

so $A^- + C - A^-ACAA^-$ is a generalized inverse of A regardless of the choice of A^- and C . Now let B be any generalized inverse of A and define $C = B - A^-$, where A^- is some particular generalized inverse of A . Then, because $ABA = A$, we have

$$\begin{aligned} A^- + C - A^-ACAA^- &= A^- + (B - A^-) - A^-A(B - A^-)AA^- \\ &= B - A^-ABAA^- + A^-AA^-AA^- \\ &= B - A^-AA^- + A^-AA^- = B, \end{aligned}$$

and so the proof is complete. \square

From Definition 5.2, we know that AA^+ is the matrix that projects vectors orthogonally onto $R(A)$, whereas A^+A projects vectors orthogonally onto $R(A')$. Theorem 5.25 looks at the matrices AA^- and A^-A .

Theorem 5.25 Let A be an $m \times n$ matrix, B be an $m \times p$ matrix, and C be a $q \times n$ matrix. Then

- (a) $AA^-B = B$ if and only if $R(B) \subset R(A)$,
- (b) $CA^-A = C$ if and only if $R(C') \subset R(A')$.

Proof. Clearly if $AA^-B = B$, the columns of B are linear combinations of the columns of A from which it follows that $R(B) \subset R(A)$. Conversely, if $R(B) \subset R(A)$ then an $n \times p$ matrix D exists, such that $B = AD$. Using this and the identity $AA^-A = A$, we have

$$B = AD = AA^-AD = AA^-B,$$

and so we have proven (a). Part (b) is proven in a similar fashion. \square

Like AA^+ , the matrix AA^- projects vectors onto the vector space $R(A)$. For instance, from Theorem 5.25, we see that $x = AA^-x$ for any $x \in R(A)$, and if $x \notin R(A)$, then $y = AA^-x \in R(A)$ because y clearly is a linear combination of the columns of A . However, although AA^+ projects vectors orthogonally onto $R(A)$, the projections given by AA^- are oblique projections unless AA^- is symmetric. In particular, AA^- projects onto $R(A)$ along the vector space $R(I_m - AA^-)$. Similarly, $(A^-A)' = A'A^{-'}$ is the projection matrix for $R(A')$ along $R(I_n - A'A^{-'})$.

We will find the following result useful in a later chapter.

Theorem 5.26 Let A , B , and C be matrices of sizes $p \times m$, $m \times n$, and $n \times q$, respectively. If $\text{rank}(ABC) = \text{rank}(B)$, then $C(ABC)^-A$ is a generalized inverse of B .

Proof. Our proof follows that of Srivastava and Khatri (1979). Using Theorem 2.8, we have

$$\text{rank}(B) = \text{rank}(ABC) \leq \text{rank}(AB) \leq \text{rank}(B)$$

and

$$\text{rank}(B) = \text{rank}(ABC) \leq \text{rank}(BC) \leq \text{rank}(B),$$

so that evidently

$$\text{rank}(AB) = \text{rank}(BC) = \text{rank}(B) = \text{rank}(ABC). \quad (5.35)$$

Using Theorem 2.10 along with the identity

$$A(BC)\{I_q - (ABC)^-ABC\} = (0),$$

we find that

$$\text{rank}(ABC) + \text{rank}(BC\{I_q - (ABC)^-ABC\}) - \text{rank}(BC) \leq \text{rank}\{(0)\} = 0,$$

so that

$$\text{rank}(BC\{I_q - (ABC)^-ABC\}) \leq \text{rank}(BC) - \text{rank}(ABC) = 0,$$

where the equality follows from (5.35). But this can be true only if

$$BC\{I_q - (ABC)^-ABC\} = \{I_q - BC(ABC)^-A\}B(C) = (0).$$

Again applying Theorem 2.10, this time on the middle expression above, we obtain

$$\text{rank}(\{I_q - BC(ABC)^-A\}B) + \text{rank}(BC) - \text{rank}(B) \leq \text{rank}\{(0)\} = 0,$$

or equivalently,

$$\text{rank}(\{I_q - BC(ABC)^-A\}B) \leq \text{rank}(B) - \text{rank}(BC) = 0,$$

where, again, the equality follows from (5.35). This implies that

$$\{I_q - BC(ABC)^-A\}B = B - B\{C(ABC)^-A\}B = (0),$$

and so the result follows. □

Some properties of a generalized inverse of $A'A$ are given in Theorem 5.27.

Theorem 5.27 Let $(A'A)^-$ be any generalized inverse of $A'A$, where A is an $m \times n$ matrix. Then

- (a) $(A'A)^{-'}$ is a generalized inverse of $A'A$,
- (b) the matrix $A(A'A)^-A'$ does not depend on the choice of the generalized inverse $(A'A)^-$,
- (c) $A(A'A)^-A'$ is symmetric even if $(A'A)^-$ is not symmetric.

Proof. Transposing the equation $A'A(A'A)^-A'A = A'A$ yields

$$A'A(A'A)^{-'}A'A = A'A,$$

so (a) follows. To prove (b) and (c), first note that

$$\begin{aligned} A(A'A)^-A'A &= AA^+A(A'A)^-A'A = (AA^+)'A(A'A)^-A'A \\ &= A^{+'}A'A(A'A)^-A'A = A^{+'}A'A \\ &= (AA^+)'A = AA^+A = A. \end{aligned}$$

Then

$$\begin{aligned} A(A'A)^-A' &= A(A'A)^-A'A^{+'}A' = A(A'A)^-A'(AA^+)' \\ &= A(A'A)^-A'AA^+ = AA^+, \end{aligned} \tag{5.36}$$

where the last equality applies the identity, $A(A'A)^-A'A = A$, just proven; (b) follows from (5.36) because A^+ , and hence also AA^+ , is unique. The symmetry of $A(A'A)^-A'$ follows from the symmetry of AA^+ . \square

We have seen that if S is the vector space spanned by the columns of the $m \times r$ matrix X_1 , then its projection matrix is given by

$$P_S = X_1(X_1'X_1)^+X_1'. \tag{5.37}$$

An immediate consequence of Theorem 5.27 is that the Moore–Penrose inverse in (5.37) can be replaced by any generalized inverse of $X_1'X_1$; that is, regardless of the choice of $(X_1'X_1)^-$, we have

$$P_S = X_1(X_1'X_1)^-X_1'.$$

We will see in Chapter 6 that the $\{1, 3\}$ -inverse is useful in finding least squares solutions to an inconsistent system of linear equations. Consequently, this inverse is commonly called the least squares inverse. We will denote a $\{1, 3\}$ -inverse of A by A^L ; that is, $A^{(1,3)} = A^L$. Since a least squares inverse of A is also a $\{1\}$ -inverse of A ,

the properties given in Theorem 5.23 also apply to A^L . Some additional properties of least squares inverses are given below.

Theorem 5.28 Let A be an $m \times n$ matrix. Then

- (a) for any least squares inverse, A^L , of A , $AA^L = AA^+$,
- (b) $(A'A)^- A'$ is a least squares inverse of A for any generalized inverse, $(A'A)^-$, of $A'A$.

Proof. Since $AA^L A = A$ and $(AA^L)' = AA^L$, we find that

$$\begin{aligned} AA^L &= AA^+ AA^L = (AA^+)'(AA^L)' = A^{+'} A' A^{L'} A' \\ &= A^{+'} (AA^L A)' = A^{+'} A' = (AA^+)' = AA^+, \end{aligned}$$

and so (a) holds. Part (b) follows from the proof of Theorem 5.27 because it was shown that

$$A(A'A)^- A' A = A,$$

which yields the first condition of a least squares inverse, and it was also shown that

$$A(A'A)^- A' = AA^+,$$

so that the symmetry of $A(A'A)^- A'$ follows from the symmetry of AA^+ . □

5.9 COMPUTING GENERALIZED INVERSES

In this section, we review some computational formulas for generalized inverses. The emphasis here is not on the development of formulas best suited for the numerical computation of generalized inverses on a computer. For instance, the most common method of computing the Moore–Penrose inverse of a matrix is through the computation of its singular value decomposition; that is, if $A = P_1 \Delta Q_1'$ is the singular value decomposition of A as given in Corollary 4.1.1, then A^+ can be easily computed via the formula $A^+ = Q_1 \Delta^{-1} P_1'$. The formulas provided here and in the problems are ones that, in some cases, may be useful for the computation of the generalized inverse of matrices of small size but, in most cases, are primarily useful for theoretical purposes.

Greville (1960) obtained an expression for the Moore–Penrose inverse of a matrix partitioned in the form $[B \quad c]$, where, of course, the matrix B and the vector c have the same number of rows. This formula can be then used recursively to compute the Moore–Penrose inverse of an $m \times n$ matrix A . To see this, let a_j denote the j th

column of A and define $A_j = (\mathbf{a}_1, \dots, \mathbf{a}_j)$, so that A_j is the $m \times j$ matrix containing the first j columns of A . Greville has shown that if we write $A_j = [A_{j-1} \quad \mathbf{a}_j]$, then

$$A_j^+ = \begin{bmatrix} A_{j-1}^+ - \mathbf{d}_j \mathbf{b}_j' \\ \mathbf{b}_j' \end{bmatrix}, \quad (5.38)$$

where $\mathbf{d}_j = A_{j-1}^+ \mathbf{a}_j$,

$$\mathbf{b}_j' = \begin{cases} (\mathbf{c}_j' \mathbf{c}_j)^{-1} \mathbf{c}_j', & \text{if } \mathbf{c}_j \neq \mathbf{0}, \\ (1 + \mathbf{d}_j' \mathbf{d}_j)^{-1} \mathbf{d}_j' A_{j-1}^+, & \text{if } \mathbf{c}_j = \mathbf{0}, \end{cases}$$

and $\mathbf{c}_j = \mathbf{a}_j - A_{j-1} \mathbf{d}_j$. Thus, $A^+ = A_n^+$ can be computed by successively computing $A_2^+, A_3^+, \dots, A_n^+$.

Example 5.10 We will use the procedure above to compute the Moore–Penrose inverse of the matrix

$$A = \begin{bmatrix} 1 & 1 & 2 & 3 \\ 1 & -1 & 0 & 1 \\ 1 & 1 & 2 & 3 \end{bmatrix}.$$

We begin by computing the inverse of $A_2 = [\mathbf{a}_1 \quad \mathbf{a}_2] = [A_1 \quad \mathbf{a}_2]$. We find that

$$\begin{aligned} A_1^+ &= (\mathbf{a}_1' \mathbf{a}_1)^{-1} \mathbf{a}_1' = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \end{bmatrix}, \\ \mathbf{d}_2 &= A_1^+ \mathbf{a}_2 = \frac{1}{3}, \\ \mathbf{c}_2 &= \mathbf{a}_2 - A_1 \mathbf{d}_2 = \mathbf{a}_2 - \frac{1}{3} \mathbf{a}_1 = \frac{1}{3} \begin{bmatrix} 2 \\ -4 \\ 2 \end{bmatrix}. \end{aligned}$$

Since $\mathbf{c}_2 \neq \mathbf{0}$, we get

$$\mathbf{b}_2' = \mathbf{c}_2^+ = (\mathbf{c}_2' \mathbf{c}_2)^{-1} \mathbf{c}_2' = \frac{1}{4} \begin{bmatrix} 1 & -2 & 1 \end{bmatrix},$$

and thus,

$$A_2^+ = \begin{bmatrix} A_1^+ - \mathbf{d}_2 \mathbf{b}_2' \\ \mathbf{b}_2' \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 1 & -2 & 1 \end{bmatrix}.$$

The inverse of $A_3 = [A_2 \quad \mathbf{a}_3]$ now can be computed by using A_2^+ and

$$\mathbf{d}_3 = A_2^+ \mathbf{a}_3 = \begin{bmatrix} 1 \\ 1 \end{bmatrix},$$

$$\mathbf{c}_3 = \mathbf{a}_3 - A_2 \mathbf{d}_3 = \begin{bmatrix} 2 \\ 0 \\ 2 \end{bmatrix} - \begin{bmatrix} 2 \\ 0 \\ 2 \end{bmatrix} = \mathbf{0}.$$

Since $\mathbf{c}_3 = \mathbf{0}$, we find that

$$\begin{aligned} \mathbf{b}'_3 &= (1 + \mathbf{d}'_3 \mathbf{d}_3)^{-1} \mathbf{d}'_3 A_2^+ = (1 + 2)^{-1} [1 \quad 1] \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 1 & -2 & 1 \end{bmatrix} \\ &= \frac{1}{6} [1 \quad 0 \quad 1], \end{aligned}$$

and so

$$A_3^+ = \begin{bmatrix} A_2^+ - \mathbf{d}_3 \mathbf{b}'_3 \\ \mathbf{b}'_3 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 1 & 6 & 1 \\ 1 & -6 & 1 \\ 2 & 0 & 2 \end{bmatrix}.$$

Finally, to obtain the Moore–Penrose inverse of $A = A_4$, we compute

$$\mathbf{d}_4 = A_3^+ \mathbf{a}_4 = \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix},$$

$$\mathbf{c}_4 = \mathbf{a}_4 - A_3 \mathbf{d}_4 = \begin{bmatrix} 3 \\ 1 \\ 3 \end{bmatrix} - \begin{bmatrix} 3 \\ 1 \\ 3 \end{bmatrix} = \mathbf{0},$$

$$\mathbf{b}'_4 = (1 + \mathbf{d}'_4 \mathbf{d}_4)^{-1} \mathbf{d}'_4 A_3^+ = \frac{1}{12} [1 \quad 2 \quad 1].$$

Consequently, the Moore–Penrose inverse of A is given by

$$A_4^+ = \begin{bmatrix} A_3^+ - \mathbf{d}_4 \mathbf{b}'_4 \\ \mathbf{b}'_4 \end{bmatrix} = \frac{1}{12} \begin{bmatrix} 0 & 4 & 0 \\ 1 & -6 & 1 \\ 1 & -2 & 1 \\ 1 & 2 & 1 \end{bmatrix}.$$

A common method of computing a generalized inverse, that is, a $\{1\}$ -inverse, of a matrix is based on the row reduction of that matrix to Hermite form.

Definition 5.4 An $m \times m$ matrix H is said to be in Hermite form if the following four conditions hold:

- (a) H is an upper triangular matrix.
- (b) h_{ii} equals 0 or 1 for each i .
- (c) If $h_{ii} = 0$, then $h_{ij} = 0$ for all j .
- (d) If $h_{ii} = 1$, then $h_{ji} = 0$ for all $j \neq i$

Before applying this concept of Hermite forms to find a generalized inverse of a matrix, we will need a couple of results regarding matrices in Hermite form. The first of these two results says that any square matrix can be transformed to a matrix in Hermite form through its premultiplication by a nonsingular matrix. Details of the proof are given in Rao (1973).

Theorem 5.29 Let A be an $m \times m$ matrix. Then a nonsingular $m \times m$ matrix C exists, such that $CA = H$, where H is in Hermite form.

The proof of Theorem 5.30 will be left to the reader as an exercise.

Theorem 5.30 Suppose the $m \times m$ matrix H is in Hermite form. Then H is idempotent; that is, $H^2 = H$.

The connection between a generalized inverse of a square matrix A and matrices in Hermite form is established in the following theorem. This result says that any matrix C satisfying the conditions of Theorem 5.29 will be a generalized inverse of A .

Theorem 5.31 Let A be an $m \times m$ matrix and C be an $m \times m$ nonsingular matrix for which $CA = H$, where H is a matrix in Hermite form. Then the matrix C is a generalized inverse of A .

Proof. We need to show that $ACA = A$. Now from Theorem 5.30, we know that H is idempotent, and so

$$CAC A = H^2 = H = CA.$$

The result then follows by premultiplying this equation by C^{-1} . □

The matrix C can be obtained by transforming A , through elementary row transformations, to a matrix in Hermite form. This process is illustrated in the following example.

Example 5.11 We will find a generalized inverse of the 3×3 matrix

$$A = \begin{bmatrix} 2 & 2 & 4 \\ 4 & -2 & 2 \\ 2 & -4 & -2 \end{bmatrix}.$$

First, we perform row transformations on A so that the resulting matrix has its first diagonal element equal to one, whereas the remaining elements in the first column are all equal to zero. This can be achieved via the matrix equation $C_1 A = A_1$, where

$$C_1 = \begin{bmatrix} 1/2 & 0 & 0 \\ -2 & 1 & 0 \\ -1 & 0 & 1 \end{bmatrix}, \quad A_1 = \begin{bmatrix} 1 & 1 & 2 \\ 0 & -6 & -6 \\ 0 & -6 & -6 \end{bmatrix}.$$

Next, we use row transformations on A_1 so that the resulting matrix has its second diagonal element equal to one, whereas each of the remaining elements in the second column is zero. This can be written as $C_2A_1 = A_2$, where

$$C_2 = \begin{bmatrix} 1 & 1/6 & 0 \\ 0 & -1/6 & 0 \\ 0 & -1 & 1 \end{bmatrix}, \quad A_2 = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix}.$$

The matrix A_2 satisfies the conditions of Definition 5.4, and so it is in Hermite form. Thus, we have $C_2A_1 = C_2C_1A = A_2$, so by Theorem 5.31 a generalized inverse of A is given by

$$C = C_2C_1 = \frac{1}{6} \begin{bmatrix} 1 & 1 & 0 \\ 2 & -1 & 0 \\ 6 & -6 & 6 \end{bmatrix}.$$

Not only is a generalized inverse not necessarily unique, but this particular method of producing a generalized inverse does not, in general, yield a unique matrix. For instance, in the second transformation given above, $C_2A_1 = A_2$, we could have chosen

$$C_2 = \begin{bmatrix} 1 & 0 & 1/6 \\ 0 & -1/6 & 0 \\ 0 & -2 & 2 \end{bmatrix}.$$

In this case, we would have obtained the generalized inverse

$$C = C_2C_1 = \frac{1}{6} \begin{bmatrix} 2 & 0 & 1 \\ 2 & -1 & 0 \\ 12 & -12 & 12 \end{bmatrix}.$$

The method of finding a generalized inverse of a matrix by transforming it to a matrix in Hermite form can be easily extended from square matrices to rectangular matrices. Theorem 5.32 indicates how such an extension is possible.

Theorem 5.32 Let A be an $m \times n$ matrix, where $m < n$. Define the matrix A_* as

$$A_* = \begin{bmatrix} A \\ (0) \end{bmatrix},$$

so that A_* is $n \times n$, and let C be any $n \times n$ nonsingular matrix for which CA_* is in Hermite form. If we partition C as $C = [C_1 \ C_2]$, where C_1 is $n \times m$, then C_1 is a generalized inverse of A .

Proof. We know from Theorem 5.31 that C is a generalized inverse of A_* . Hence, $A_*CA_* = A_*$. Simplifying the left-hand side of this identity, we find that

$$\begin{aligned} A_*CA_* &= \begin{bmatrix} A \\ (0) \end{bmatrix} [C_1 \quad C_2] \begin{bmatrix} A \\ (0) \end{bmatrix} \\ &= \begin{bmatrix} AC_1 & AC_2 \\ (0) & (0) \end{bmatrix} \begin{bmatrix} A \\ (0) \end{bmatrix} \\ &= \begin{bmatrix} AC_1A \\ (0) \end{bmatrix}. \end{aligned}$$

Equating this identity to A_* , we get $AC_1A = A$, and so the proof is complete. \square

Clearly, an analogous result holds for the case in which $m > n$.

Example 5.12 Suppose that we wish to find a generalized inverse of the matrix

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 0 & 1 \\ 1 & 1 & 2 \\ 2 & 0 & 2 \end{bmatrix}.$$

Consequently, we consider the augmented matrix

$$A_* = [A \quad \mathbf{0}] = \begin{bmatrix} 1 & 1 & 2 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 2 & 0 \\ 2 & 0 & 2 & 0 \end{bmatrix}.$$

Proceeding as in the previous example, we obtain a nonsingular matrix C so that CA_* is in Hermite form. One such matrix is given by

$$C = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ -1 & 0 & 1 & 0 \\ 0 & -2 & 0 & 1 \end{bmatrix}.$$

Thus, partitioning this matrix as

$$C = \begin{bmatrix} C_1 \\ C_2 \end{bmatrix},$$

we find that a generalized inverse of A is given by

$$C_1 = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & -1 & 0 & 0 \\ -1 & 0 & 1 & 0 \end{bmatrix}.$$

A least squares generalized inverse of a matrix A can be computed by first computing a generalized inverse of $A'A$ and then using the relationship, $A^L = (A'A)^- A'$, established in Theorem 5.28(b).

Example 5.13 To find a least squares inverse of the matrix A from Example 5.12, we first compute

$$A'A = \begin{bmatrix} 7 & 2 & 9 \\ 2 & 2 & 4 \\ 9 & 4 & 13 \end{bmatrix}.$$

By transforming this matrix to Hermite form, we find that a generalized inverse of $A'A$ is given by

$$(A'A)^- = \frac{1}{10} \begin{bmatrix} 2 & -2 & 0 \\ -2 & 7 & 0 \\ -10 & -10 & 10 \end{bmatrix}.$$

Hence, a least squares inverse of A is given by

$$A^L = (A'A)^- A' = \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

PROBLEMS

5.1 Prove results (a)–(d) of Theorem 5.3.

5.2 Use Theorem 5.3(h) to find the Moore–Penrose inverse of

$$A = \begin{bmatrix} 1 & 1 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 1 \\ 2 & 0 & 1 \end{bmatrix}.$$

5.3 Find the Moore–Penrose inverse of the vector

$$\mathbf{a} = \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}.$$

5.4 Provide the proofs for (f)–(j) of Theorem 5.3.

5.5 Prove Theorem 5.6.

5.6 Use the spectral decomposition of the matrix

$$A = \begin{bmatrix} 2 & 0 & 1 \\ 0 & 2 & 3 \\ 1 & 3 & 5 \end{bmatrix}$$

to find its Moore–Penrose inverse.

5.7 Consider the matrix

$$A = \begin{bmatrix} 0 & -1 & 2 \\ 0 & -1 & 2 \\ 3 & 2 & -1 \end{bmatrix}.$$

- (a) Find the Moore–Penrose inverse of AA' , and then use Theorem 5.3(g) to find A^+ .
 - (b) Use A^+ to find the projection matrix for the range of A and the projection matrix for the row space of A .
- 5.8** Show that the converse of Theorem 5.5(c) does not hold. That is, give an example of a symmetric matrix A for which $A^+ = A$ yet A is not idempotent.
- 5.9** Let A be an $m \times n$ matrix with $\text{rank}(A) = 1$. Show that $A^+ = c^{-1}A'$, where $c = \text{tr}(A'A)$.
- 5.10** Let \mathbf{x} and \mathbf{y} be $m \times 1$ vectors, and let $\mathbf{1}_m$ be the $m \times 1$ vector with each element equal to one. Obtain expressions for the Moore–Penrose inverses of
- (a) $\mathbf{1}_m \mathbf{1}_m'$,
 - (b) $I_m - m^{-1} \mathbf{1}_m \mathbf{1}_m'$,
 - (c) $\mathbf{x} \mathbf{x}'$,
 - (d) $\mathbf{x} \mathbf{y}'$.
- 5.11** Let A be an $m \times n$ matrix. Show that each of the matrices, AA^+ , A^+A , $(I_m - AA^+)$, and $(I_n - A^+A)$ is idempotent.
- 5.12** Let A be an $m \times n$ matrix. Establish the following identities:
- (a) $A'AA^+ = A^+AA' = A'$.
 - (b) $A'A^+A^+ = A^+A^+A' = A^+$.
 - (c) $A(A'A)^+A'A = AA'(AA')^+A = A$.
- 5.13** Let A be an $m \times n$ matrix and B be an $n \times n$ positive definite matrix. Show that

$$ABA'(ABA')^+A = A.$$

5.14 Let A be an $m \times n$ matrix. Show that

- (a) $AB = (0)$ if and only if $B^+A^+ = (0)$, where B is an $n \times p$ matrix.
 - (b) $A^+B = (0)$ if and only if $A'B = (0)$, where B is an $m \times p$ matrix.
- 5.15** Let A be an $m \times m$ symmetric matrix having rank r . Show that if A has one nonzero eigenvalue λ of multiplicity r , then $A^+ = \lambda^{-2}A$.

- 5.16** Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that if B has full row rank, then

$$AB(AB)^+ = AA^+.$$

- 5.17** Let A be an $m \times m$ symmetric matrix. Show that

- (a) if A is nonnegative definite, then so is A^+ ,
 (b) if $Ax = 0$ for some vector x , then $A^+x = 0$ also.

- 5.18** Let A be an $m \times m$ symmetric matrix with $\text{rank}(A) = r$. Use the spectral decomposition of A to show that if B is any $m \times m$ symmetric matrix with $\text{rank}(B) = m - r$, such that $AB = (0)$, then $A^+A + B^+B = I_m$.

- 5.19** Let A and B be $m \times m$ nonnegative definite matrices, and suppose that $A - B$ is also nonnegative definite. Show that $B^+ - A^+$ is nonnegative definite if and only if $\text{rank}(A) = \text{rank}(B)$.

- 5.20** Let A be an $m \times n$ matrix and B be an $n \times m$ matrix. Suppose that $\text{rank}(A) = \text{rank}(B)$ and, further, that the space spanned by the eigenvectors corresponding to the positive eigenvalues of $A'A$ is the same as that spanned by the eigenvectors corresponding to the positive eigenvalues of BB' . Show that $(AB)^+ = B^+A^+$.

- 5.21** Prove Theorem 5.8.

- 5.22** Prove (b)–(d) of Theorem 5.10.

- 5.23** For each case below use Theorem 5.10 to determine whether $(AB)^+ = B^+A^+$.

$$(a) \quad A = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix}.$$

$$(b) \quad A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}.$$

- 5.24** Let A be an $m \times n$ matrix and B be an $n \times m$ matrix. Show that $(AB)^+ = B^+A^+$ if $A'ABB' = BB'A'A$.

- 5.25** Prove Theorem 5.14.

- 5.26** Find the Moore–Penrose inverse of the matrix

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

- 5.27** Use Corollary 5.13.1(d) to find the Moore–Penrose inverse of the matrix $A = [U \quad V]$, where

$$U = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \quad V = \begin{bmatrix} 1 & -2 \\ -1 & 1 \\ 0 & 1 \end{bmatrix}.$$

- 5.28** Use Corollary 5.13.1(c) to find the Moore–Penrose inverse of the matrix $A = [U \ V]$, where

$$U = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} 2 & 2 \\ 2 & 0 \\ -1 & 0 \\ 1 & -2 \\ 0 & 1 \end{bmatrix}.$$

- 5.29** Let the vectors w , x , y , and z be given by

$$w = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad x = \begin{bmatrix} 1 \\ 1 \\ -2 \\ 0 \end{bmatrix}, \quad y = \begin{bmatrix} 1 \\ -1 \\ 0 \\ 0 \end{bmatrix}, \quad z = \begin{bmatrix} 1 \\ 1 \\ 1 \\ -3 \end{bmatrix}.$$

Use Theorem 5.17 to find the Moore–Penrose inverse of the matrix $A = wx' + yz'$.

- 5.30** Suppose A is an $m \times n$ matrix while c and d are $m \times 1$ and $n \times 1$ vectors. Let $u = (I_m - AA^+)c$ and $v = d'(I_n - A^+A)$.

(a) Show that

$$(A + cd')^+ = A^+ - A^+cu^+ - v^+d'A^+ + (1 + d'A^+c)v^+u^+,$$

if $u \neq 0$ and $v \neq 0$.

(b) Show that

$$(A + cd')^+ = A^+ - A^+c(A^+c)^+A^+ - v^+d'A^+,$$

if $u = 0$, $v \neq 0$, and $1 + d'A^+c = 0$.

(c) Show that

$$(A + cd')^+ = A^+ - (1 + d'A^+c)^{-1}A^+cd'A^+,$$

if $u = 0$, $v = 0$, and $1 + d'A^+c \neq 0$.

- 5.31** If the $m \times 1$ random vector x has a multinomial distribution (see, for example, Johnson, et al. 1997), then $\text{var}(x_i) = np_i(1 - p_i)$ and $\text{cov}(x_i, x_j) = -np_ip_j$ for $i \neq j$, where n is a positive integer and $0 < p_i < 1$, such that $p_1 + \cdots + p_m = 1$. If Ω is the covariance matrix of x , use Theorem 5.18 to show that Ω is singular and

$$\Omega^+ = n^{-1}(I_m - m^{-1}\mathbf{1}_m\mathbf{1}_m')D^{-1}(I_m - m^{-1}\mathbf{1}_m\mathbf{1}_m'),$$

where $D = \text{diag}(p_1, \dots, p_m)$.

- 5.32** Let A be an $m \times m$ nonsingular symmetric matrix and \mathbf{c} and \mathbf{d} be $m \times 1$ vectors. Show that if $A + \mathbf{cd}'$ is singular, then it has A^{-1} as a generalized inverse.
- 5.33** Let A be an $m \times m$ symmetric matrix, and suppose that \mathbf{c} and \mathbf{d} are $m \times 1$ vectors that are in the column space of A . Show that if $1 + \mathbf{d}'A^+\mathbf{c} \neq 0$

$$(A + \mathbf{cd}')^+ = A^+ - \frac{A^+\mathbf{cd}'A^+}{1 + \mathbf{d}'A^+\mathbf{c}}.$$

- 5.34** Find a generalized inverse, different than the Moore–Penrose inverse, of the vector given in Problem 5.3.
- 5.35** Consider the diagonal matrix $A = \text{diag}(0, 2, 3)$.
- (a) Find a generalized inverse of A having rank of 2.
 - (b) Find a generalized inverse of A that has rank of 3 and is diagonal.
 - (c) Find a generalized inverse of A that is not diagonal.
- 5.36** Let A be an $m \times m$ matrix partitioned as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

where A_{11} is $r \times r$. Show that if $\text{rank}(A) = \text{rank}(A_{11}) = r$, then

$$\begin{bmatrix} A_{11}^{-1} & (0) \\ (0) & (0) \end{bmatrix}$$

is a generalized inverse of A .

- 5.37** Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that B^-A^- will be a generalized inverse of AB for any choice of A^- and B^- if $\text{rank}(B) = n$.
- 5.38** Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that for any choice of A^- and B^- , B^-A^- will be a generalized inverse of AB if and only if A^-ABB^- is idempotent.
- 5.39** Show that a matrix B is a generalized inverse of A if and only if AB is idempotent and $\text{rank}(A) = \text{rank}(AB)$.
- 5.40** Let A , P , and Q be $m \times n$, $p \times m$, and $n \times q$ matrices, respectively. Show that if P has full column rank and Q has full row rank, then $Q^-A^-P^-$ is a generalized inverse of PAQ .
- 5.41** Let A be an $m \times n$ matrix, B be an $m \times m$ matrix, and C be an $n \times n$ matrix. Show that if B and C are nonsingular, then an $n \times m$ matrix D is a generalized inverse of BAC if and only if $D = C^{-1}A^-B^{-1}$ for some generalized inverse A^- of A .
- 5.42** Show that the matrix AA^- yields orthogonal projections if and only if it is symmetric; that is, show that

$$(\mathbf{x} - AA^-\mathbf{x})'AA^-\mathbf{x} = 0$$

for all \mathbf{x} if and only if AA^- is symmetric.

5.43 A matrix B is called a reflexive generalized inverse of A if it satisfies the first two conditions of a Moore–Penrose inverse; that is, B is a reflexive inverse of A if $ABA = A$ and $BAB = B$. Show the following.

- (a) A generalized inverse B of a matrix A is reflexive if and only if $\text{rank}(B) = \text{rank}(A)$.
 (b) For any two matrices E and F ,

$$B = Q \begin{bmatrix} \Delta^{-1} & E \\ F & F\Delta E \end{bmatrix} P'$$

is a reflexive inverse of A , where A has the singular value decomposition given by

$$A = P \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} Q'.$$

5.44 Suppose that the $m \times n$ matrix A is partitioned as $A = [A_1 \ A_2]$, where A_1 is $m \times r$, and $\text{rank}(A) = \text{rank}(A_1) = r$. Show that $A(A'A)^-A' = A_1(A_1'A_1)^-A_1'$.

5.45 Suppose that the $m \times n$ matrix A has been partitioned as $A' = [U' \ V']$, where U is $m_1 \times n$, V is $m_2 \times n$ and $m_1 + m_2 = m$. Show that a generalized inverse of A is given by $A^- = [W \ X]$, where

$$W = (I_n - (I_n - U^-U)\{V(I_n - U^-U)\}^-V)U^-, \\ X = (I_n - U^-U)\{V(I_n - U^-U)\}^-.$$

5.46 Use the recursive procedure described in Section 9 to obtain the Moore–Penrose inverse of the matrix.

$$A = \begin{bmatrix} 1 & -1 & -1 \\ -1 & 1 & 1 \\ 2 & -1 & 1 \end{bmatrix}.$$

5.47 Find a generalized inverse of the matrix A in the previous problem by finding a nonsingular matrix that transforms it into a matrix having Hermite form.

5.48 Find a generalized inverse of the matrix

$$A = \begin{bmatrix} 1 & -1 & -2 & 1 \\ -2 & 4 & 3 & -2 \\ 1 & 1 & -3 & 1 \end{bmatrix}.$$

5.49 Find a least squares inverse for the matrix A given in the previous problem.

5.50 Let A be an $m \times m$ matrix and C be an $m \times m$ nonsingular matrix such that $CA = H$ is in Hermite form. Show that A is idempotent if and only if H is a generalized inverse of A .

- 5.51** Let A be an $m \times n$ matrix and suppose that B is an $n \times n$ matrix that satisfies conditions 1 and 4 of the Moore–Penrose inverse of $A'A$. Show that $A^+ = BA'$.
- 5.52** Let A be an $m \times n$ matrix and B be an $n \times m$ matrix. Show that B is the Moore–Penrose inverse of A if and only if B is a least squares inverse of A and A is a least squares inverse of B .
- 5.53** Let A be an $m \times n$ matrix. Show that for any least squares generalized inverse, $(AA')^L$, of AA' , $A^+ = A'(AA')^L$.
- 5.54** Let A be an $m \times n$ matrix. Show that an $n \times m$ matrix B is a least squares generalized inverse of A if and only if $A'AB = A'$.
- 5.55** Let A be an $m \times n$ matrix, and let $(AA')^-$ and $(A'A)^-$ be any generalized inverses of AA' and $A'A$, respectively. Show that

$$A^+ = A'(AA')^-A(A'A)^-A'.$$

- 5.56** It was shown in Theorem 5.31 that a generalized inverse of an $m \times m$ matrix A can be obtained by finding a nonsingular matrix that row reduces A to Hermite form. Show that there is a similar result for column reduction to Hermite form; that is, show that if C is a nonsingular matrix, such that $AC = H$, where H is in Hermite form, then C is a generalized inverse of A .
- 5.57** Prove Theorem 5.30.
- 5.58** Let A be an $m \times n$ matrix. Show that

$$A^+ = \lim_{\delta \rightarrow 0} (A'A + \delta^2 I_n)^{-1} A' = \lim_{\delta \rightarrow 0} A'(AA' + \delta^2 I_m)^{-1}.$$

- 5.59** Penrose (1956) obtained the following recursive method for calculating the Moore–Penrose generalized inverse of an $m \times n$ matrix A . Successively calculate B_2, B_3, \dots , where

$$B_{i+1} = i^{-1} \text{tr}(B_i A' A) I_n - B_i A' A$$

and B_1 is defined to be the $n \times n$ identity matrix. If $\text{rank}(A) = r$, then $B_{r+1} A' A = (0)$ and

$$A^+ = r \{ \text{tr}(B_r A' A) \}^{-1} B_r A'.$$

Use this method to compute the Moore–Penrose inverse of the matrix A of Example 5.10.

- 5.60** Let λ be the largest eigenvalue of AA' , where A is an $m \times n$ matrix. Let α be any constant satisfying $0 < \alpha < 2/\lambda$, and define $X_1 = \alpha A'$. Ben-Israel (1966) has shown that if we define

$$X_{i+1} = X_i(2I_m - AX_i)$$

for $i = 1, 2, \dots$, then $X_i \rightarrow A^+$ as $i \rightarrow \infty$. Use this iterative procedure to compute the Moore–Penrose inverse of the matrix A of Example 5.10 on a computer. Stop the iterative process when

$$\text{tr}\{(X_{i+1} - X_i)'(X_{i+1} - X_i)\}$$

gets small. Note that λ does not need to be computed because we must have

$$\frac{2}{\text{tr}(AA')} < \frac{2}{\lambda}.$$

- 5.61** Use the results of Section 5 to obtain the expression given in (5.38) for the Moore–Penrose inverse of the matrix $A_j = [A_{j-1} \quad \mathbf{a}_j]$.

6

SYSTEMS OF LINEAR EQUATIONS

6.1 INTRODUCTION

As mentioned at the beginning of Chapter 5, one of the applications of generalized inverses is in finding solutions to a system of linear equations of the form

$$A\mathbf{x} = \mathbf{c}, \quad (6.1)$$

where A is an $m \times n$ matrix of constants, \mathbf{c} is an $m \times 1$ vector of constants, and \mathbf{x} is an $n \times 1$ vector of variables for which solutions are needed. There are three possibilities regarding the solution \mathbf{x} to (6.1): there is no solution, there is exactly one solution, there is more than one solution. We will see in Section 6.3 that if there is more than one solution, then there are actually infinitely many solutions. In this chapter, we discuss such issues as the existence of solutions to (6.1), the form of a general solution, and the number of linearly independent solutions. We also look at the special application of finding least squares solutions to (6.1), when an exact solution does not exist.

6.2 CONSISTENCY OF A SYSTEM OF EQUATIONS

In this section, we will obtain necessary and sufficient conditions for the existence of a vector \mathbf{x} satisfying (6.1). When one or more such vectors exist, the system of

equations is said to be consistent; otherwise, the system is referred to as an inconsistent system. Our first necessary and sufficient condition for consistency is that the vector \mathbf{c} is in the column space of A or, equivalently, that the rank of the augmented matrix $[A \ \mathbf{c}]$ is the same as the rank of A .

Theorem 6.1 The system of equations, $A\mathbf{x} = \mathbf{c}$, is consistent if and only if $\text{rank}([A \ \mathbf{c}]) = \text{rank}(A)$.

Proof. If $\mathbf{a}_1, \dots, \mathbf{a}_n$ are the columns of A , then the equation $A\mathbf{x} = \mathbf{c}$ can be written as

$$A\mathbf{x} = [\mathbf{a}_1 \ \cdots \ \mathbf{a}_n] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n x_i \mathbf{a}_i = \mathbf{c}.$$

Clearly, this equation holds for some \mathbf{x} if and only if \mathbf{c} is a linear combination of the columns of A , in which case $\text{rank}[A \ \mathbf{c}] = \text{rank}(A)$. \square

Example 6.1 Consider the system of equations that has

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix}.$$

Clearly, the rank of A is 2, whereas

$$|[A \ \mathbf{c}]| = \begin{vmatrix} 1 & 2 & 1 \\ 2 & 1 & 5 \\ 1 & 0 & 3 \end{vmatrix} = 0,$$

so that the rank of $[A \ \mathbf{c}]$ is also 2. Thus, we know from Theorem 6.1 that the system of equations $A\mathbf{x} = \mathbf{c}$ is consistent.

Although Theorem 6.1 is useful in determining whether a given system of linear equations is consistent, it does not tell us how to find a solution to the system when it is consistent. Theorem 6.2 gives an alternative necessary and sufficient condition for consistency applying a generalized inverse, A^- , of A . An obvious consequence of this result is that when the system $A\mathbf{x} = \mathbf{c}$ is consistent, then a solution will be given by $\mathbf{x} = A^-\mathbf{c}$.

Theorem 6.2 The system of equations $A\mathbf{x} = \mathbf{c}$ is consistent if and only if for some generalized inverse, A^- , of A , $AA^-\mathbf{c} = \mathbf{c}$.

Proof. First, suppose that the system is consistent and \mathbf{x}_* is a solution, so that $\mathbf{c} = A\mathbf{x}_*$. Premultiplying this identity by AA^- , where A^- is any generalized inverse of A , yields

$$AA^-\mathbf{c} = AA^-A\mathbf{x}_* = A\mathbf{x}_* = \mathbf{c},$$

as is required. Conversely, now suppose that there is a generalized inverse of A satisfying $AA^-c = c$. Define $x_* = A^-c$, and note that

$$Ax_* = AA^-c = c.$$

Thus, because $x_* = A^-c$ is a solution, the system is consistent, and so the proof is complete. \square

Suppose that A_1 and A_2 are any two generalized inverses of A so that $AA_1A = AA_2A = A$. In addition, suppose that A_1 satisfies the condition of Theorem 6.2; that is, $AA_1c = c$. Then A_2 satisfies the same condition because

$$AA_2c = AA_2(AA_1c) = (AA_2A)A_1c = AA_1c = c.$$

Thus, in applying Theorem 6.2, one will need to check the given condition for only one generalized inverse of A , and it does not matter which generalized inverse is used. In particular, we can use the Moore–Penrose inverse, A^+ , of A .

Corollary 6.2.1 and Corollary 6.2.2 involve some special cases regarding the matrix A .

Corollary 6.2.1 If A is an $m \times m$ nonsingular matrix and c is an $m \times 1$ vector of constants, then the system $Ax = c$ is consistent.

Corollary 6.2.2 If the $m \times n$ matrix A has rank equal to m , then the system $Ax = c$ is consistent.

Proof. Since A has full row rank, it follows from Theorem 5.23(f) that $AA^- = I_m$. As a result, $AA^-c = c$, and so from Theorem 6.2, the system must be consistent. \square

Example 6.2 Consider the system of equations $Ax = c$, where

$$A = \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix}, \quad c = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}.$$

A generalized inverse of the transpose of A was given in Example 5.12. Using this inverse, we find that

$$\begin{aligned} AA^-c &= \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}. \end{aligned}$$

Since this is \mathbf{c} , the system of equations is consistent, and a solution is given by

$$A^{-}\mathbf{c} = \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix}.$$

The generalized inverse given in Example 5.12 is not the only generalized inverse for A so we could have solved this problem using a different generalized inverse. For instance, it is easily verified that

$$A^{-} = \begin{bmatrix} 3 & 3 & -1 \\ 1 & -1 & 0 \\ -4 & -3 & 2 \\ 0 & 0 & 0 \end{bmatrix}$$

satisfies $AA^{-}A = A$. Using this choice for A^{-} , we again find that the consistency condition holds because

$$\begin{aligned} AA^{-}\mathbf{c} &= \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 3 & 3 & -1 \\ 1 & -1 & 0 \\ -4 & -3 & 2 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \\ &= \begin{bmatrix} 0 & -1 & 1 \\ -1 & 0 & 1 \\ -1 & -1 & 2 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}. \end{aligned}$$

However, we get a different solution because

$$A^{-}\mathbf{c} = \begin{bmatrix} 3 & 3 & -1 \\ 1 & -1 & 0 \\ -4 & -3 & 2 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 10 \\ 1 \\ -8 \\ 0 \end{bmatrix}.$$

The system of linear equations $A\mathbf{x} = \mathbf{c}$ is a special case of the more general system of linear equations given by $AXB = C$, where A is $m \times n$, B is $p \times q$, C is $m \times q$, and X is $n \times p$. A necessary and sufficient condition for the existence of a solution matrix X satisfying this system is given in Theorem 6.3.

Theorem 6.3 Let A , B , and C be matrices of constants, where A is $m \times n$, B is $p \times q$, and C is $m \times q$. Then the system of equations,

$$AXB = C,$$

is consistent if and only if for some generalized inverses A^- and B^- ,

$$AA^-CB^-B = C. \quad (6.2)$$

Proof. Suppose that the system is consistent and the matrix X_* is a solution, so that $C = AX_*B$. Premultiplying by AA^- and postmultiplying by B^-B , where A^- and B^- are any generalized inverses of A and B , we find that

$$AA^-CB^-B = AA^-AX_*BB^-B = AX_*B = C,$$

and so (6.2) holds. On the other hand, if A^- and B^- satisfy (6.2), define $X_* = A^-CB^-$, and note that X_* is a solution because

$$AX_*B = AA^-CB^-B = C,$$

so the proof is complete. \square

Using an argument similar to that given after Theorem 6.2, we can verify that if (6.2) is satisfied for any one particular choice of A^- and B^- , then it will hold for all choices of A^- and B^- . Consequently, the application of Theorem 6.3 is not dependent on the choices of generalized inverses for A and B .

6.3 SOLUTIONS TO A CONSISTENT SYSTEM OF EQUATIONS

We have seen that if the system of equations $Ax = c$ is consistent, then $x = A^-c$ is a solution regardless of the choice of the generalized inverse A^- . Thus, if A^-c is not the same for all choices of A^- , then our system of equations has more than one solution. In fact, we will see that even when A^-c does not depend on the choice of A^- , which is the case if $c = 0$, our system of equations may have many solutions. Theorem 6.4 gives a general expression for all solutions to the system.

Theorem 6.4 Suppose that $Ax = c$ is a consistent system of equations, and let A^- be any generalized inverse of the $m \times n$ matrix A . Then, for any $n \times 1$ vector y ,

$$x_y = A^-c + (I_n - A^-A)y \quad (6.3)$$

is a solution, and for any solution, x_* , a vector y exists, such that $x_* = x_y$.

Proof. Since $Ax = c$ is a consistent system of equations, we know from Theorem 6.2 that $AA^-c = c$, and so

$$\begin{aligned} Ax_y &= AA^-c + A(I_n - A^-A)y \\ &= c + (A - AA^-A)y = c, \end{aligned}$$

because $AA^-A = A$. Thus, \mathbf{x}_y is a solution regardless of the choice of \mathbf{y} . On the other hand, if \mathbf{x}_* is an arbitrary solution, so that $A\mathbf{x}_* = \mathbf{c}$, it follows that $A^-A\mathbf{x}_* = A^-\mathbf{c}$. Consequently,

$$A^-\mathbf{c} + (I_n - A^-A)\mathbf{x}_* = A^-\mathbf{c} + \mathbf{x}_* - A^-A\mathbf{x}_* = \mathbf{x}_*,$$

so that $\mathbf{x}_* = \mathbf{x}_{x_*}$. This completes the proof. \square

The set of solutions given in Theorem 6.4 is expressed in terms of a fixed generalized inverse A^- and an arbitrary $n \times 1$ vector \mathbf{y} . Alternatively, this set of all solutions can be expressed in terms of an arbitrary generalized inverse of A .

Corollary 6.4.1 Suppose that $A\mathbf{x} = \mathbf{c}$ is a consistent system of equations, where $\mathbf{c} \neq \mathbf{0}$. If B is a generalized inverse of A , then $\mathbf{x} = B\mathbf{c}$ is a solution, and for any solution \mathbf{x}_* , a generalized inverse B exists, such that $\mathbf{x}_* = B\mathbf{c}$.

Proof. Theorem 6.4 was not dependent on the choice of the generalized inverse, so by choosing $A^- = B$ and $\mathbf{y} = \mathbf{0}$ in (6.3), we prove that $\mathbf{x} = B\mathbf{c}$ is a solution. All that remains to be shown is that for any particular A^- and \mathbf{y} , we can find a generalized inverse B , such that the expression in (6.3) equals $B\mathbf{c}$. Now because $\mathbf{c} \neq \mathbf{0}$, it has at least one component, say c_i , not equal to 0. Define the $n \times m$ matrix C as $C = c_i^{-1}\mathbf{y}\mathbf{e}_i'$, so that $C\mathbf{c} = \mathbf{y}$. Since the system of equations $A\mathbf{x} = \mathbf{c}$ is consistent, we must have $AA^-\mathbf{c} = \mathbf{c}$, and so

$$\begin{aligned} \mathbf{x}_y &= A^-\mathbf{c} + (I_n - A^-A)\mathbf{y} = A^-\mathbf{c} + (I_n - A^-A)C\mathbf{c} \\ &= A^-\mathbf{c} + C\mathbf{c} - A^-AC\mathbf{c} = A^-\mathbf{c} + C\mathbf{c} - A^-ACAA^-\mathbf{c} \\ &= (A^- + C - A^-ACAA^-)\mathbf{c}. \end{aligned}$$

However, it follows from Theorem 5.24 that $A^- + C - A^-ACAA^-$ is a generalized inverse of A for any choice of the $n \times m$ matrix C , and so the proof is complete. \square

Example 6.3 For the consistent system of equations discussed in Example 6.2, we have

$$\begin{aligned} A^-A &= \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix} \\ &= \begin{bmatrix} -1 & -1 & -1 & -2 \\ 0 & 1 & 0 & 2 \\ 2 & 1 & 2 & 2 \\ 0 & 0 & 0 & 0 \end{bmatrix}, \end{aligned}$$

where we have used the first of the two generalized inverses given in that example. Consequently, a general solution to this system of equations is given by

$$\begin{aligned} \mathbf{x}_y &= A^- \mathbf{c} + (I_4 - A^- A) \mathbf{y} \\ &= \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 & 1 & 1 & 2 \\ 0 & 0 & 0 & -2 \\ -2 & -1 & -1 & -2 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ &= \begin{bmatrix} -3 + 2y_1 + y_2 + y_3 + 2y_4 \\ 1 - 2y_4 \\ 5 - 2y_1 - y_2 - y_3 - 2y_4 \\ y_4 \end{bmatrix}, \end{aligned}$$

where \mathbf{y} is an arbitrary 4×1 vector.

Our next theorem gives a result, analogous to Theorem 6.4, for the system of equations $AXB = C$. The proof will be left to the reader as an exercise.

Theorem 6.5 Let $AXB = C$ be a consistent system of equations, where A is $m \times n$, B is $p \times q$, and C is $m \times q$. Then for any generalized inverses, A^- and B^- , and any $n \times p$ matrix, Y ,

$$X_Y = A^- C B^- + Y - A^- A Y B B^-$$

is a solution, and for any solution, X_* , a matrix Y exists, such that $X_* = X_Y$.

Example 6.4 Suppose an $m \times m$ nonsingular matrix C is partitioned as $C = [A \ B]$, where A is $m \times m_1$ and B is $m \times m_2$. We wish to find an expression for C^{-1} in terms of the submatrices A and B . Partition C^{-1} as $C^{-1} = [X' \ Y']'$, where X is $m_1 \times m$ and Y is $m_2 \times m$. The system of equations $XA = I_{m_1}$ is consistent since $A^- A = I_{m_1}$, and so it follows from Theorem 6.5 that

$$X = A^- + U(I_{m_1} - AA^-), \quad (6.4)$$

where U is an arbitrary $m_1 \times m$ matrix. Note that Theorem 5.23(e) and $(I_{m_1} - AA^-)A = (0)$ guarantee that $\text{rank}(I_{m_1} - AA^-) = m_2$, and since the columns of B are linearly independent of those in A , we must also have $\text{rank}\{(I_{m_1} - AA^-)B\} = m_2$. Thus, $XB = A^- B + U(I_{m_1} - AA^-)B = (0)$ is consistent, and another application of Theorem 6.5 yields

$$\begin{aligned} U &= -A^- B \{(I_{m_1} - AA^-)B\}^- \\ &\quad + V(I_{m_1} - \{(I_{m_1} - AA^-)B\} \{(I_{m_1} - AA^-)B\}^-), \end{aligned} \quad (6.5)$$

where V is an arbitrary $m_1 \times m$ matrix. Since $(I_m - AA^-)B$ and $(I_m - AA^-)$ have the same rank, they also have the same column space, so it follows from Theorem 5.25 that

$$\{(I_m - AA^-)B\}\{(I_m - AA^-)B\}^-(I_m - AA^-) = (I_m - AA^-).$$

Using this while substituting (6.5) into (6.4), we get

$$X = A^- - A^-B\{(I_m - AA^-)B\}^-(I_m - AA^-).$$

In a similar fashion, using the equations $YB = I_{m_2}$ and $YA = (0)$, we can show that

$$Y = B^- - B^-A\{(I_m - BB^-)A\}^-(I_m - BB^-).$$

In some applications, it may be important to know whether a consistent system of equations yields a unique solution; that is, under what conditions will (6.3) yield the same solution for all choices of \mathbf{y} ?

Theorem 6.6 If $A\mathbf{x} = \mathbf{c}$ is a consistent system of equations, then the solution $\mathbf{x}_* = A^-\mathbf{c}$ is a unique solution if and only if $A^-A = I_n$, where A^- is any generalized inverse of the $m \times n$ matrix A .

Proof. Note that $\mathbf{x}_* = A^-\mathbf{c}$ is a unique solution if and only if $\mathbf{x}_\mathbf{y} = \mathbf{x}_*$ for all choices of \mathbf{y} , where $\mathbf{x}_\mathbf{y}$ is as defined in (6.3). In other words, the solution is unique if and only if

$$(I_n - A^-A)\mathbf{y} = \mathbf{0}$$

for all \mathbf{y} , and clearly this is equivalent to the condition $(I_n - A^-A) = (0)$ or $A^-A = I_n$. \square

We saw in Theorem 5.23(g) that $\text{rank}(A) = n$ if and only if $A^-A = I_n$. As a result, we can restate the necessary and sufficient condition of Theorem 6.6 as follows in Corollary 6.6.1.

Corollary 6.6.1 Suppose that $A\mathbf{x} = \mathbf{c}$ is a consistent system of equations. Then the solution $\mathbf{x}_* = A^-\mathbf{c}$ is a unique solution if and only if $\text{rank}(A) = n$.

Example 6.5 We saw in Example 6.1 that the system of equations $A\mathbf{x} = \mathbf{c}$, where

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix},$$

is consistent. The Moore–Penrose inverse of the transpose of A was obtained in Example 5.1. Using this inverse, we find that

$$\begin{aligned} A^+A &= \frac{1}{14} \begin{bmatrix} -3 & 6 & 5 \\ 8 & -2 & -4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix} \\ &= \frac{1}{14} \begin{bmatrix} 14 & 0 \\ 0 & 14 \end{bmatrix} = I_2. \end{aligned}$$

Thus, the system of equations $A\mathbf{x} = \mathbf{c}$ has the unique solution given by

$$\begin{aligned} A^+\mathbf{c} &= \frac{1}{14} \begin{bmatrix} -3 & 6 & 5 \\ 8 & -2 & -4 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix} \\ &= \frac{1}{14} \begin{bmatrix} 42 \\ -14 \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \end{bmatrix}. \end{aligned}$$

Suppose that a system of linear equations has more than one solution, and let \mathbf{x}_1 and \mathbf{x}_2 be two different solutions. Then, because $A\mathbf{x}_i = \mathbf{c}$ for $i = 1$ and 2 , it follows that for any scalar α

$$A\{\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2\} = \alpha A\mathbf{x}_1 + (1 - \alpha)A\mathbf{x}_2 = \alpha\mathbf{c} + (1 - \alpha)\mathbf{c} = \mathbf{c}.$$

Thus, $\mathbf{x} = \{\alpha\mathbf{x}_1 + (1 - \alpha)\mathbf{x}_2\}$ is also a solution. Since α was arbitrary, we see that if a system has more than one solution, then it has infinitely many solutions. However, the number of linearly independent solutions to a consistent system of equations having $\mathbf{c} \neq \mathbf{0}$ must be between 1 and n ; that is, a set of linearly independent solutions $\{\mathbf{x}_1, \dots, \mathbf{x}_r\}$ exists, such that every solution can be expressed as a linear combination of the solutions, $\mathbf{x}_1, \dots, \mathbf{x}_r$. In other words, any solution \mathbf{x} can be written as $\mathbf{x} = \alpha_1\mathbf{x}_1 + \dots + \alpha_r\mathbf{x}_r$, for some coefficients $\alpha_1, \dots, \alpha_r$. Note that because $A\mathbf{x}_i = \mathbf{c}$ for each i , we must have

$$A\mathbf{x} = A\left(\sum_{i=1}^r \alpha_i\mathbf{x}_i\right) = \sum_{i=1}^r \alpha_i A\mathbf{x}_i = \sum_{i=1}^r \alpha_i\mathbf{c} = \left(\sum_{i=1}^r \alpha_i\right)\mathbf{c},$$

and so if \mathbf{x} is a solution, the coefficients must satisfy the identity $\alpha_1 + \dots + \alpha_r = 1$. Theorem 6.7 tells us exactly how to determine this number of linearly independent solutions r when $\mathbf{c} \neq \mathbf{0}$. We will delay the discussion of the situation in which $\mathbf{c} = \mathbf{0}$ until the next section.

Theorem 6.7 Suppose that the system $A\mathbf{x} = \mathbf{c}$ is consistent, where A is $m \times n$ and $\mathbf{c} \neq \mathbf{0}$. Then each solution can be expressed as a linear combination of r linearly independent solutions, where $r = n - \text{rank}(A) + 1$.

Proof. Using (6.3) with the particular generalized inverse A^+ , we begin with the $n + 1$ solutions, $\mathbf{x}_0 = A^+\mathbf{c}$, $\mathbf{x}_{e_1} = A^+\mathbf{c} + (I_n - A^+A)\mathbf{e}_1, \dots, \mathbf{x}_{e_n} = A^+\mathbf{c} + (I_n - A^+A)\mathbf{e}_n$, where, as usual, \mathbf{e}_i denotes the $n \times 1$ vector whose only nonzero element is 1 in the i th position. Now every solution can be expressed as a linear combination of these solutions because for any $\mathbf{y} = (y_1, \dots, y_n)'$,

$$\mathbf{x}_y = A^+\mathbf{c} + (I_n - A^+A)\mathbf{y} = \left(1 - \sum_{i=1}^n y_i\right) \mathbf{x}_0 + \sum_{i=1}^n y_i \mathbf{x}_{e_i}.$$

Thus, if we define the $n \times (n + 1)$ matrix $X = (\mathbf{x}_0, \mathbf{x}_{e_1}, \dots, \mathbf{x}_{e_n})$, the proof will be complete if we can show that $\text{rank}(X) = n - \text{rank}(A) + 1$. Note that we can write X as $X = BC$, where B and C are the $n \times (n + 1)$ and $(n + 1) \times (n + 1)$ matrices given by $B = (A^+\mathbf{c}, I_n - A^+A)$ and

$$C = \begin{bmatrix} 1 & \mathbf{1}'_n \\ \mathbf{0} & I_n \end{bmatrix}.$$

Clearly, C is nonsingular because it is upper triangular and the product of its diagonal elements is 1. Consequently, from Theorem 1.10, we know that $\text{rank}(X) = \text{rank}(B)$. Note also that

$$\begin{aligned} (I_n - A^+A)'A^+\mathbf{c} &= (I_n - A^+A)A^+\mathbf{c} = (A^+ - A^+AA^+)\mathbf{c} \\ &= (A^+ - A^+)\mathbf{c} = \mathbf{0}, \end{aligned}$$

so that the first column of B is orthogonal to the remaining columns. This implies that

$$\text{rank}(B) = \text{rank}(A^+\mathbf{c}) + \text{rank}(I_n - A^+A) = 1 + \text{rank}(I_n - A^+A),$$

because the consistency condition $AA^+\mathbf{c} = \mathbf{c}$ and $\mathbf{c} \neq \mathbf{0}$ guarantee that $A^+\mathbf{c} \neq \mathbf{0}$. All that remains is to show that $\text{rank}(I_n - A^+A) = n - \text{rank}(A)$. Now because A^+A is the projection matrix of $R(A^+) = R(A')$, it follows that $I_n - A^+A$ is the projection matrix of the orthogonal complement of $R(A')$ or, in other words, the null space of A , $N(A)$. Since $\dim\{N(A)\} = n - \text{rank}(A)$, we must have $\text{rank}(I_n - A^+A) = n - \text{rank}(A)$. \square

Since $\mathbf{x}_0 = A^+\mathbf{c}$ is orthogonal to the columns of $(I_n - A^+A)$, when constructing a set of r linearly independent solutions, one of these solutions always will be \mathbf{x}_0 , with the remaining solutions given by \mathbf{x}_y for $r - 1$ different choices of $\mathbf{y} \neq \mathbf{0}$. This statement is not dependent on the choice of A^+ as the generalized inverse in (6.3), because $A^-\mathbf{c}$ and $(I_n - A^-A)\mathbf{y}$ are linearly independent regardless of the choice of A^- if $\mathbf{c} \neq \mathbf{0}$, $\mathbf{y} \neq \mathbf{0}$. The proof of this linear independence is left as an exercise.

Example 6.6 We saw that the system of equations $A\mathbf{x} = \mathbf{c}$ of Examples 6.2 and 6.3 has the set of solutions consisting of all vectors of the form

$$\mathbf{x}_{\mathbf{y}} = A^{-}\mathbf{c} + (I_4 - A^{-}A)\mathbf{y} = \begin{bmatrix} -3 + 2y_1 + y_2 + y_3 + 2y_4 \\ 1 - 2y_4 \\ 5 - 2y_1 - y_2 - y_3 - 2y_4 \\ y_4 \end{bmatrix}.$$

Since the last row of the 3×4 matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix}$$

is the sum of the first two rows, $\text{rank}(A) = 2$. Thus, the system of equations possesses

$$n - \text{rank}(A) + 1 = 4 - 2 + 1 = 3$$

linearly independent solutions. Three linearly independent solutions can be obtained through appropriate choices of the \mathbf{y} vector. For instance, because $A^{-}\mathbf{c}$ and $(I_4 - A^{-}A)\mathbf{y}$ are linearly independent, the three solutions

$$A^{-}\mathbf{c}, \quad A^{-}\mathbf{c} + (I_4 - A^{-}A)_{\cdot i}, \quad A^{-}\mathbf{c} + (I_4 - A^{-}A)_{\cdot j}$$

will be linearly independent if the i th and j th columns of $(I_4 - A^{-}A)$ are linearly independent. Looking back at the matrix $(I_4 - A^{-}A)$ given in Example 6.3, we see that its first and fourth columns are linearly independent. Thus, three linearly independent solutions of $A\mathbf{x} = \mathbf{c}$ are given by

$$\begin{aligned} A^{-}\mathbf{c} &= \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix}, \\ A^{-}\mathbf{c} + (I_4 - A^{-}A)_{\cdot 1} &= \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \\ -2 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 3 \\ 0 \end{bmatrix}, \\ A^{-}\mathbf{c} + (I_4 - A^{-}A)_{\cdot 4} &= \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ -2 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ 3 \\ 1 \end{bmatrix}. \end{aligned}$$

6.4 HOMOGENEOUS SYSTEMS OF EQUATIONS

The system of equations $Ax = c$ is called a nonhomogeneous system of equations when $c \neq \mathbf{0}$, whereas $Ax = \mathbf{0}$ is referred to as a homogeneous system of equations. In this section, we obtain some results regarding homogeneous systems of equations. One obvious distinction between homogeneous and nonhomogeneous systems is that a homogeneous system of equations must be consistent because it will always have the trivial solution, $x = \mathbf{0}$. A homogeneous system will then have a unique solution only when the trivial solution is the only solution. Conditions for the existence of nontrivial solutions, which we state in the next theorem, follow directly from Theorem 6.6 and Corollary 6.6.1.

Theorem 6.8 Suppose that A is an $m \times n$ matrix. The system $Ax = \mathbf{0}$ has nontrivial solutions if and only if $A^-A \neq I_n$, or equivalently if and only if $\text{rank}(A) < n$.

If the system $Ax = \mathbf{0}$ has more than one solution, and $\{x_1, \dots, x_r\}$ is a set of r solutions, then $x = \alpha_1 x_1 + \dots + \alpha_r x_r$ is also a solution regardless of the choice of $\alpha_1, \dots, \alpha_r$, because

$$Ax = A \left(\sum_{i=1}^r \alpha_i x_i \right) = \sum_{i=1}^r \alpha_i Ax_i = \sum_{i=1}^r \alpha_i \mathbf{0} = \mathbf{0}.$$

In fact, we have the following.

Theorem 6.9 If A is an $m \times n$ matrix, then the set of all solutions to the system of equations $Ax = \mathbf{0}$ forms a vector subspace of R^n having dimension $n - \text{rank}(A)$.

Proof. The result follows immediately from the fact that the set of all solutions of $Ax = \mathbf{0}$ is the null space of A . \square

In contrast to Theorem 6.9, the set of all solutions to a nonhomogeneous system of equations will not form a vector subspace. This is because, as we have seen in the previous section, a linear combination of solutions to a nonhomogeneous system yields another solution only if the coefficients sum to one. Additionally, a nonhomogeneous system cannot have $\mathbf{0}$ as a solution.

The general form of a solution given in Theorem 6.4 applies to both homogeneous and nonhomogeneous systems. Thus, for any $n \times 1$ vector y ,

$$x_y = (I_n - A^-A)y$$

is a solution to the system $Ax = \mathbf{0}$, and for any solution, x_* , a vector y exists, such that $x_* = x_y$. Theorem 6.10 shows that the set of solutions of $Ax = c$ can be expressed in terms of the set of solutions to $Ax = \mathbf{0}$.

Theorem 6.10 Let \mathbf{x}_* be any solution to the system of equations $A\mathbf{x} = \mathbf{c}$. Then

- (a) if $\mathbf{x}_\#$ is a solution to the system $A\mathbf{x} = \mathbf{0}$, $\mathbf{x} = \mathbf{x}_* + \mathbf{x}_\#$ is a solution of $A\mathbf{x} = \mathbf{c}$, and
- (b) for any solution \mathbf{x} to the equation $A\mathbf{x} = \mathbf{c}$, a solution $\mathbf{x}_\#$ to the equation $A\mathbf{x} = \mathbf{0}$ exists, such that $\mathbf{x} = \mathbf{x}_* + \mathbf{x}_\#$.

Proof. Note that if $\mathbf{x}_\#$ is as defined in (a), then

$$A(\mathbf{x}_* + \mathbf{x}_\#) = A\mathbf{x}_* + A\mathbf{x}_\# = \mathbf{c} + \mathbf{0} = \mathbf{c},$$

and so $\mathbf{x} = \mathbf{x}_* + \mathbf{x}_\#$ is a solution to $A\mathbf{x} = \mathbf{c}$. To prove (b), define $\mathbf{x}_\# = \mathbf{x} - \mathbf{x}_*$, so that $\mathbf{x} = \mathbf{x}_* + \mathbf{x}_\#$. Then because $A\mathbf{x} = \mathbf{c}$ and $A\mathbf{x}_* = \mathbf{c}$, it follows that

$$A\mathbf{x}_\# = A(\mathbf{x} - \mathbf{x}_*) = A\mathbf{x} - A\mathbf{x}_* = \mathbf{c} - \mathbf{c} = \mathbf{0},$$

and so the proof is complete. \square

Our next result, regarding the number of linearly independent solutions possessed by a homogeneous system of equations, follows immediately from Theorem 6.9.

Theorem 6.11 Each solution of the homogeneous system of equations $A\mathbf{x} = \mathbf{0}$ can be expressed as a linear combination of r linearly independent solutions, where $r = n - \text{rank}(A)$.

Example 6.7 Consider the system of equations $A\mathbf{x} = \mathbf{0}$, where

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix}.$$

We saw in Example 6.5 that $A^+A = I_2$. Thus, the system only has the trivial solution $\mathbf{0}$.

Example 6.8 Since the matrix

$$A = \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix}$$

from Example 6.6 has rank of 2, the homogeneous system of equations $A\mathbf{x} = \mathbf{0}$ has $r = n - \text{rank}(A) = 4 - 2 = 2$ linearly independent solutions. Any set of two

linearly independent columns of the matrix $(I_4 - A^-A)$ will be a set of linearly independent solutions; for example, the first and fourth columns,

$$\begin{bmatrix} 2 \\ 0 \\ -2 \\ 0 \end{bmatrix}, \quad \begin{bmatrix} 2 \\ -2 \\ -2 \\ 1 \end{bmatrix},$$

are linearly independent solutions.

6.5 LEAST SQUARES SOLUTIONS TO A SYSTEM OF LINEAR EQUATIONS

In some situations in which we have an inconsistent system of equations $Ax = c$, it may be desirable to find the vector or set of vectors that comes “closest” to satisfying the equations. If x_* is one choice for x , then x_* will approximately satisfy our system of equations if $Ax_* - c$ is close to 0 . One of the most common ways of measuring the closeness of $Ax_* - c$ to 0 is through the computation of the sum of squares of the components of the vector $Ax_* - c$. Any vector minimizing this sum of squares is referred to as a least squares solution.

Definition 6.1 The $n \times 1$ vector x_* is said to be a least squares solution to the system of equations $Ax = c$ if the inequality

$$(Ax_* - c)'(Ax_* - c) \leq (Ax - c)'(Ax - c) \quad (6.6)$$

holds for every $n \times 1$ vector x .

Of course, we have already used the concept of a least squares solution in many of our examples on regression analysis. In particular, we have seen that if the matrix X has full column rank, then the least squares solution for $\hat{\beta}$ in the fitted regression equation, $\hat{y} = X\hat{\beta}$, is given by $\hat{\beta} = (X'X)^{-1}X'y$. The generalized inverses that we have discussed in Chapter 5 will enable us to obtain a unified treatment of this problem including cases in which X is not of full rank.

In Section 5.8, we briefly discussed the $\{1, 3\}$ -inverse of a matrix A , that is, any matrix satisfying the first and third conditions of the Moore–Penrose inverse. We referred to this type of inverse as a least squares inverse of A . Theorem 6.12 motivates this description.

Theorem 6.12 Let A^L be any $\{1, 3\}$ -inverse of a matrix A . Then the vector $x_* = A^L c$ is a least squares solution to the system of equations $Ax = c$.

Proof. We must show that (6.6) holds when $\mathbf{x}_* = A^L \mathbf{c}$. The right-hand side of (6.6) can be written as

$$\begin{aligned}
 (A\mathbf{x} - \mathbf{c})'(A\mathbf{x} - \mathbf{c}) &= \{(A\mathbf{x} - AA^L \mathbf{c}) + (AA^L \mathbf{c} - \mathbf{c})\}' \\
 &\quad \times \{(A\mathbf{x} - AA^L \mathbf{c}) + (AA^L \mathbf{c} - \mathbf{c})\} \\
 &= (A\mathbf{x} - AA^L \mathbf{c})'(A\mathbf{x} - AA^L \mathbf{c}) \\
 &\quad + (AA^L \mathbf{c} - \mathbf{c})'(AA^L \mathbf{c} - \mathbf{c}) \\
 &\quad + 2(A\mathbf{x} - AA^L \mathbf{c})'(AA^L \mathbf{c} - \mathbf{c}) \\
 &\geq (AA^L \mathbf{c} - \mathbf{c})'(AA^L \mathbf{c} - \mathbf{c}) \\
 &= (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}),
 \end{aligned}$$

where the inequality follows from the fact that

$$(A\mathbf{x} - AA^L \mathbf{c})'(A\mathbf{x} - AA^L \mathbf{c}) \geq 0$$

and

$$\begin{aligned}
 (A\mathbf{x} - AA^L \mathbf{c})'(AA^L \mathbf{c} - \mathbf{c}) &= (\mathbf{x} - A^L \mathbf{c})' A' (AA^L \mathbf{c} - \mathbf{c}) \\
 &= (\mathbf{x} - A^L \mathbf{c})' A' ((AA^L)' \mathbf{c} - \mathbf{c}) \\
 &= (\mathbf{x} - A^L \mathbf{c})' (A' A^{L'} A' \mathbf{c} - A' \mathbf{c}) \\
 &= (\mathbf{x} - A^L \mathbf{c})' (A' \mathbf{c} - A' \mathbf{c}) = 0.
 \end{aligned} \tag{6.7}$$

This completes the proof. \square

Corollary 6.12.1 The vector \mathbf{x}_* is a least squares solution to the system $A\mathbf{x} = \mathbf{c}$ if and only if

$$(A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) = \mathbf{c}'(I_m - AA^L)\mathbf{c}.$$

Proof. From Theorem 6.12, $A^L \mathbf{c}$ is a least squares solution for any choice of A^L , and its sum of squared errors is given by

$$\begin{aligned}
 (AA^L \mathbf{c} - \mathbf{c})'(AA^L \mathbf{c} - \mathbf{c}) &= \mathbf{c}'(AA^L - I_m)'(AA^L - I_m)\mathbf{c} \\
 &= \mathbf{c}'(AA^L - I_m)^2 \mathbf{c} \\
 &= \mathbf{c}'(AA^L AA^L - 2AA^L + I_m)\mathbf{c} \\
 &= \mathbf{c}'(AA^L - 2AA^L + I_m)\mathbf{c} \\
 &= \mathbf{c}'(I_m - AA^L)\mathbf{c}.
 \end{aligned}$$

The result now follows because, by definition, a least squares solution minimizes the sum of squared errors, and so any other vector \mathbf{x}_* will be a least squares solution if and only if its sum of squared errors is equal to this minimum sum of squares, $\mathbf{c}'(I_m - AA^L)\mathbf{c}$. \square

Example 6.9 Let the system of equations $A\mathbf{x} = \mathbf{c}$ have A and \mathbf{c} given by

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 0 & 1 \\ 1 & 1 & 2 \\ 2 & 0 & 2 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix}.$$

In Example 5.13, we computed the least squares inverse,

$$A^L = \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}.$$

Since

$$AA^L\mathbf{c} = \frac{1}{10} \begin{bmatrix} 5 & 0 & 5 & 0 \\ 0 & 2 & 0 & 4 \\ 5 & 0 & 5 & 0 \\ 0 & 4 & 0 & 8 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 5 \\ 2.2 \\ 5 \\ 4.4 \end{bmatrix} \neq \mathbf{c},$$

it follows from Theorem 6.2 that the system of equations is inconsistent. A least squares solution is then given by

$$A^L\mathbf{c} = \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 2.2 \\ 2.8 \\ 0 \end{bmatrix}.$$

Since $(AA^L\mathbf{c} - \mathbf{c})' = (5, 2.2, 5, 4.4)' - (4, 1, 6, 5)' = (1, 1.2, -1, -0.6)'$, the sum of squared errors for the least squares solution is

$$(AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}) = 3.8.$$

In general, a least squares solution is not unique. For instance, the reader can easily verify that the matrix

$$B = \begin{bmatrix} -2 & -0.8 & -2 & -1.6 \\ -1.5 & -1.2 & -1.5 & -2.4 \\ 2 & 1 & 2 & 2 \end{bmatrix}$$

is also a least squares inverse of A . Consequently,

$$Bc = \begin{bmatrix} -2 & -0.8 & -2 & -1.6 \\ -1.5 & -1.2 & -1.5 & -2.4 \\ 2 & 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} = \begin{bmatrix} -28.8 \\ -28.2 \\ 31 \end{bmatrix}$$

is another least squares solution. However, because $(ABc - c)' = (5, 2.2, 5, 4.4)' - (4, 1, 6, 5)' = (1, 1.2, -1, -0.6)'$, the sum of squared errors for this least squares solution is, as it must be, identical to that of the previous solution.

The following result will be useful in establishing the general form of a least squares solution. It indicates that although a least squares solution \mathbf{x}_* may not be unique, the vector $A\mathbf{x}_*$ will be unique.

Theorem 6.13 The vector \mathbf{x}_* is a least squares solution to the system $A\mathbf{x} = \mathbf{c}$ if and only if

$$A\mathbf{x}_* = AA^L\mathbf{c}. \quad (6.8)$$

Proof. Using Theorem 6.2, we see that the system of equations given in (6.8) is consistent because

$$AA^L(AA^L\mathbf{c}) = (AA^LA)A^L\mathbf{c} = AA^L\mathbf{c}.$$

The sum of squared errors for any vector \mathbf{x}_* satisfying (6.8) is

$$\begin{aligned} (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) &= (AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}) \\ &= \mathbf{c}'(AA^L - I_m)^2\mathbf{c} \\ &= \mathbf{c}'(I_m - AA^L)\mathbf{c}, \end{aligned}$$

so by Corollary 6.12.1, \mathbf{x}_* is a least squares solution. Conversely, now suppose that \mathbf{x}_* is a least squares solution. Then from Corollary 6.12.1 we must have

$$\begin{aligned} (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) &= \mathbf{c}'(I_m - AA^L)\mathbf{c} \\ &= \mathbf{c}'(I_m - AA^L)'(I_m - AA^L)\mathbf{c} \\ &= (AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}), \end{aligned} \quad (6.9)$$

where we have used the fact that $(I_m - AA^L)$ is symmetric and idempotent. However, we also have

$$\begin{aligned} (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) &= \{(A\mathbf{x}_* - AA^L\mathbf{c}) + (AA^L\mathbf{c} - \mathbf{c})\}' \\ &\quad \times \{(A\mathbf{x}_* - AA^L\mathbf{c}) + (AA^L\mathbf{c} - \mathbf{c})\} \\ &= (A\mathbf{x}_* - AA^L\mathbf{c})'(A\mathbf{x}_* - AA^L\mathbf{c}) \\ &\quad + (AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}), \end{aligned} \quad (6.10)$$

because $(A\mathbf{x}_* - AA^L\mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}) = 0$, as shown in (6.7). Now (6.9) and (6.10) imply that

$$(A\mathbf{x}_* - AA^L\mathbf{c})(A\mathbf{x}_* - AA^L\mathbf{c}) = 0,$$

which can be true only if

$$(A\mathbf{x}_* - AA^L\mathbf{c}) = 0,$$

which establishes (6.8). \square

We now give an expression in Theorem 6.14 for a general least squares solution to a system of equations.

Theorem 6.14 Let A^L be any $\{1, 3\}$ -inverse of the $m \times n$ matrix A . Define the vector

$$\mathbf{x}_y = A^L\mathbf{c} + (I_n - A^LA)\mathbf{y},$$

where \mathbf{y} is an arbitrary $n \times 1$ vector. Then, for each \mathbf{y} , \mathbf{x}_y is a least squares solution to the system of equations $A\mathbf{x} = \mathbf{c}$, and for any least squares solution \mathbf{x}_* a vector \mathbf{y} exists, such that $\mathbf{x}_* = \mathbf{x}_y$.

Proof. Since

$$A(I_n - A^LA)\mathbf{y} = (A - AA^LA)\mathbf{y} = (A - A)\mathbf{y} = \mathbf{0},$$

we have $A\mathbf{x}_y = AA^L\mathbf{c}$, and so by Theorem 6.13, \mathbf{x}_y is a least squares solution. Conversely, if \mathbf{x}_* is an arbitrary least squares solution, then by using Theorem 6.13 again, we must have

$$A\mathbf{x}_* = AA^L\mathbf{c},$$

which, when premultiplied by A^L , implies that

$$\mathbf{0} = -A^LA(\mathbf{x}_* - A^L\mathbf{c}).$$

Adding \mathbf{x}_* to both sides of this identity and then rearranging, we get

$$\begin{aligned} \mathbf{x}_* &= \mathbf{x}_* - A^LA(\mathbf{x}_* - A^L\mathbf{c}) \\ &= A^L\mathbf{c} + \mathbf{x}_* - A^L\mathbf{c} - A^LA(\mathbf{x}_* - A^L\mathbf{c}) \\ &= A^L\mathbf{c} + (I_n - A^LA)(\mathbf{x}_* - A^L\mathbf{c}). \end{aligned}$$

This completes the proof because we have shown that $\mathbf{x}_* = \mathbf{x}_y$, where $\mathbf{y} = (\mathbf{x}_* - A^L\mathbf{c})$. \square

We saw in Example 6.9 that least squares solutions are not necessarily unique. Theorem 6.14 can be used to obtain a necessary and sufficient condition for the solution to be unique.

Theorem 6.15 If A is an $m \times n$ matrix, then the system of equations $Ax = c$ has a unique least squares solution if and only if $\text{rank}(A) = n$.

Proof. It follows immediately from Theorem 6.14 that the least squares solution is unique if and only if $(I_n - A^L A) = (0)$, or equivalently, $A^L A = I_n$. The result now follows from Theorem 5.23(g). \square

In some applications, when there is more than one least squares solution to the system $Ax = c$, one may be interested in finding the least squares solution that has the smallest size. We will refer to such a solution as the minimal least squares solution. Our next result indicates that this solution is given by $x = A^+ c$.

Theorem 6.16 Suppose $x_* \neq A^+ c$ is a least squares solution to the system of equations $Ax = c$. Then

$$x'_* x_* > c' A^+ A^+ c.$$

Proof. It follows from Theorem 6.14 that x_* can be written as $x_* = A^+ c + (I_n - A^+ A)y$ for some $n \times 1$ vector y . Since $x_* \neq A^+ c$, we must have $(I_n - A^+ A)y \neq 0$, or equivalently $\{(I_n - A^+ A)y\}'(I_n - A^+ A)y = y'(I_n - A^+ A)y > 0$. The result now follows since

$$\begin{aligned} x'_* x_* &= c' A^+ A^+ c + y'(I_n - A^+ A)y + 2y'(I_n - A^+ A)A^+ c \\ &= c' A^+ A^+ c + y'(I_n - A^+ A)y \\ &> c' A^+ A^+ c. \end{aligned}$$

\square

Even when the least squares solution to a system is not unique, certain linear combinations of the elements of least squares solutions may be unique. This is the subject of our next theorem.

Theorem 6.17 Let x_* be a least squares solution to the system of equations $Ax = c$. Then $a'x_*$ is unique if and only if a is in the row space of A .

Proof. Using Theorem 6.14, if $a'x_*$ is unique regardless of the choice of the least squares solution x_* , then

$$a'x_y = a' A^L c + a'(I_n - A^L A)y$$

is the same for all choices of y . But this implies that

$$a'(I_n - A^L A) = 0'. \quad (6.11)$$

Now if (6.11) holds, then

$$a' = b' A,$$

where $\mathbf{b}' = \mathbf{a}' A^L$, and so \mathbf{a} is in the row space of A . On the other hand, if \mathbf{a} is in the row space of A , then some vector \mathbf{b} exists, such that $\mathbf{a}' = \mathbf{b}' A$. This implies that

$$\mathbf{a}'(I_n - A^L A) = \mathbf{b}' A(I_n - A^L A) = \mathbf{b}'(A - A A^L A) = \mathbf{b}'(A - A) = \mathbf{0}',$$

and so the least squares solution must be unique. \square

Example 6.10 We will obtain the general least squares solution to the system of equations presented in Example 6.9. First, note that

$$A^L A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix},$$

so that

$$\begin{aligned} \mathbf{x}_y &= A^L \mathbf{c} + (I_3 - A^L A) \mathbf{y} \\ &= \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} + \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \\ &= \begin{bmatrix} 2.2 - y_3 \\ 2.8 - y_3 \\ y_3 \end{bmatrix} \end{aligned}$$

is a least squares solution for any choice of y_3 . The quantity $\mathbf{a}' \mathbf{x}_y$ does not depend on the choice of y_3 as long as \mathbf{a} is in the row space of A ; in this case, that corresponds to \mathbf{a} being orthogonal to the vector $(-1, -1, 1)'$.

6.6 LEAST SQUARES ESTIMATION FOR LESS THAN FULL RANK MODELS

In all of our previous examples of least squares estimation for a model of the form

$$\mathbf{y} = X\beta + \epsilon, \tag{6.12}$$

where \mathbf{y} is $N \times 1$, X is $N \times m$, β is $m \times 1$, and ϵ is $N \times 1$, we have assumed that $\text{rank}(X) = m$. In this case, the normal equations,

$$X'X\hat{\beta} = X'\mathbf{y}, \tag{6.13}$$

yield a unique solution, the unique least squares estimator of β , given by

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

However, in many applications, the matrix X has less than full rank.

Example 6.11 Consider the univariate one-way classification model, which was written as

$$y_{ij} = \mu_i + \epsilon_{ij},$$

in Example 3.16, where $i = 1, \dots, k$ and $j = 1, \dots, n_i$. This model can be written in the form of (6.12), where $\beta = (\mu_1, \dots, \mu_k)'$ and

$$X = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}.$$

In this case, X is of full rank, and so

$$\hat{\beta} = (X'X)^{-1}X'\mathbf{y} = \bar{\mathbf{y}} = \left(\sum y_{1j}/n_1, \dots, \sum y_{kj}/n_k \right)'.$$

An alternative way of writing this one-way classification model is

$$y_{ij} = \mu + \tau_i + \epsilon_{ij},$$

which has $k + 1$ parameters instead of k . Here μ represents an overall effect, whereas τ_i is an effect due to treatment i . In some respects, this form of the model is more natural in that the reduced model, which has all treatment means identical, is simply a submodel with some of the parameters equal to 0, that is, $\tau_1 = \dots = \tau_k = 0$. If this second form of the model is written as $\mathbf{y} = X_*\beta_* + \epsilon$, then $\beta_* = (\mu, \tau_1, \dots, \tau_k)'$ and

$$X_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}.$$

Thus, this second parameterization of the one-way classification model has the design matrix X_* less than full rank because $\text{rank}(X_*) = k$.

In this section, we will apply some of the results of this chapter to the estimation of parameters in the model given by (6.12) when X is less than full rank. First of all, let us consider the task of solving the normal equations given by (6.13); that is, using our usual notation for a system of equations, we want to solve $A\mathbf{x} = \mathbf{c}$, where

$A = X'X$, $\mathbf{x} = \hat{\beta}$, and $\mathbf{c} = X'\mathbf{y}$. Now from Theorem 6.2, we see that (6.13) is a consistent system of equations because

$$\begin{aligned} X'X(X'X)^+X'\mathbf{y} &= X'XX^+X^+X'\mathbf{y} = X'XX^+(XX^+)'\mathbf{y} \\ &= X'XX^+XX^+\mathbf{y} = X'XX^+\mathbf{y} \\ &= X'(XX^+)'\mathbf{y} = X'X^+X'\mathbf{y} = X'\mathbf{y}. \end{aligned}$$

Consequently, using Theorem 6.4, we find that the general solution $\hat{\beta}$ can be written as

$$\hat{\beta} = (X'X)^-X'\mathbf{y} + \{I_m - (X'X)^-X'X\}\mathbf{u}, \quad (6.14)$$

or, if we use the Moore–Penrose generalized inverse, as

$$\begin{aligned} \hat{\beta} &= (X'X)^+X'\mathbf{y} + \{I_m - (X'X)^+X'X\}\mathbf{u} \\ &= X^+\mathbf{y} + (I_m - X^+X)\mathbf{u}, \end{aligned}$$

where \mathbf{u} is an arbitrary $m \times 1$ vector. The same general solution can be obtained by applying the least squares results of Section 6.5 on the system of equations

$$\mathbf{y} = X\hat{\beta}.$$

Thus, using Theorem 6.14 with $A = X$, $\mathbf{x} = \hat{\beta}$, and $\mathbf{c} = \mathbf{y}$, the least squares solution is given by

$$\hat{\beta} = X^L\mathbf{y} + (I_m - X^LX)\mathbf{u},$$

which is, of course, equivalent to that given by (6.14).

One key difference between the full rank model and the less than full rank model is that the least squares solution is unique only if X has full rank. When X is less than full rank, the model $\mathbf{y} = X\beta + \epsilon$ is overparameterized, and so not all of the parameters or linear functions of the parameters are uniquely defined; this is what leads to the infinitely many solutions for $\hat{\beta}$. Thus, when estimating linear functions of the parameters, we must make sure that we are trying to estimate a function of the parameters that is uniquely defined. This leads to the following definition of what is known as an estimable function.

Definition 6.2 The linear function $\mathbf{a}'\beta$ of the parameter vector β is estimable if and only if some $N \times 1$ vector \mathbf{b} exists, such that

$$\mathbf{a}'\beta = E(\mathbf{b}'\mathbf{y}) = \mathbf{b}'E(\mathbf{y}) = \mathbf{b}'X\beta;$$

that is, if and only if there exists a linear function of the components of \mathbf{y} , $\mathbf{b}'\mathbf{y}$, which is an unbiased estimator of $\mathbf{a}'\beta$.

The condition that a linear function $\mathbf{a}'\beta$ be estimable is equivalent to the condition that the corresponding estimator $\mathbf{a}'\hat{\beta}$ be unique. To see this, note that from Definition 6.2, the function $\mathbf{a}'\beta$ is estimable if and only if \mathbf{a} is in the row space of X , whereas it follows from Theorem 6.17 that $\mathbf{a}'\hat{\beta}$ is unique if and only if \mathbf{a} is in the row space of X . In addition, because $X'(XX')^+X$ is the projection matrix for the row space of X , we get the more practical condition for estimability of $\mathbf{a}'\beta$ given by

$$X'(XX')^+X\mathbf{a} = \mathbf{a}. \quad (6.15)$$

It follows from Theorems 5.3 and 5.28 that

$$X'(XX')^+X = X'X^{+'} = X'X^{L'} = X'(XX')^-X,$$

and so (6.15) is not dependent on the Moore–Penrose inverse as the choice of the generalized inverse of XX' .

Finally, we will demonstrate the invariance of the vector of fitted values $\hat{\mathbf{y}} = X\hat{\beta}$ and its sum of squared errors $(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}})$ to the choice of the least squares solution $\hat{\beta}$. Since $XX^+X = X$,

$$\begin{aligned} \hat{\mathbf{y}} &= X\hat{\beta} = X\{X^+\mathbf{y} + (I_m - X^+X)\mathbf{u}\} \\ &= XX^+\mathbf{y} + (X - XX^+X)\mathbf{u} = XX^+\mathbf{y}, \end{aligned}$$

which does not depend on the vector \mathbf{u} . Thus, $\hat{\mathbf{y}}$ is unique, whereas the uniqueness of

$$(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) = \mathbf{y}'(I_m - XX^+)\mathbf{y}$$

follows immediately from the uniqueness of $\hat{\mathbf{y}}$.

Example 6.12 Let us return to the one-way classification model

$$\mathbf{y} = X_*\beta_* + \epsilon$$

of Example 6.11, where $\beta_* = (\mu, \tau_1, \dots, \tau_k)'$ and

$$X_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}.$$

Since the rank of the $n \times (k+1)$ matrix X_* , where $n = \sum n_i$, is k , the least squares solution for β_* is not unique. To find the form of the general solution, note that

$$X_*'X_* = \begin{bmatrix} n & \mathbf{n}' \\ \mathbf{n} & D_n \end{bmatrix},$$

whereas a generalized inverse is given by

$$(X'_*X_*)^- = \begin{bmatrix} n^{-1} & \mathbf{0}' \\ \mathbf{0} & D_n^{-1} - n^{-1}\mathbf{1}_k\mathbf{1}_k' \end{bmatrix},$$

where $\mathbf{n} = (n_1, \dots, n_k)'$ and $D_n = \text{diag}(n_1, \dots, n_k)$. Thus, using (6.14) we have the general solution

$$\begin{aligned} \hat{\beta}_* &= \begin{bmatrix} n^{-1} & \mathbf{0}' \\ \mathbf{0} & D_n^{-1} - n^{-1}\mathbf{1}_k\mathbf{1}_k' \end{bmatrix} \begin{bmatrix} n\bar{y} \\ D_n\bar{\mathbf{y}} \end{bmatrix} \\ &\quad + \left\{ I_{k+1} - \begin{bmatrix} 1 & n^{-1}\mathbf{n}' \\ \mathbf{0} & I_k - n^{-1}\mathbf{1}_k\mathbf{n}' \end{bmatrix} \right\} \mathbf{u} \\ &= \begin{bmatrix} \bar{y} \\ \bar{\mathbf{y}} - \bar{y}\mathbf{1}_k \end{bmatrix} + \begin{bmatrix} 0 & -n^{-1}\mathbf{n}' \\ \mathbf{0} & n^{-1}\mathbf{1}_k\mathbf{n}' \end{bmatrix} \mathbf{u}, \end{aligned}$$

where $\bar{\mathbf{y}} = (\bar{y}_1, \dots, \bar{y}_k)'$ and $\bar{y} = \sum n_i \bar{y}_i / n$. Choosing $\mathbf{u} = \mathbf{0}$, we get the particular least squares solution that has $\hat{\mu} = \bar{y}$ and $\hat{\tau}_i = \bar{y}_i - \bar{y}$ for $i = 1, \dots, k$. Since $\mathbf{a}'\beta_*$ is estimable only if \mathbf{a} is in the row space of X , we find that the k quantities, $\mu + \tau_i$, $i = 1, \dots, k$, as well as any linear combinations of these quantities, are estimable. In particular, because $\mu + \tau_i = \mathbf{a}'_i\beta_*$, where $\mathbf{a}_i = (1, e'_i)'$, its estimator is given by

$$\mathbf{a}'_i\hat{\beta}_* = [1 \quad e'_i] \begin{bmatrix} \bar{y} \\ \bar{\mathbf{y}} - \bar{y}\mathbf{1}_k \end{bmatrix} = \bar{y}_i.$$

The vector of fitted values is

$$\hat{\mathbf{y}} = X_*\hat{\beta}_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix} \begin{bmatrix} \bar{y} \\ \bar{y}_1 - \bar{y} \\ \bar{y}_2 - \bar{y} \\ \vdots \\ \bar{y}_k - \bar{y} \end{bmatrix} = \begin{bmatrix} \bar{y}_1\mathbf{1}_{n_1} \\ \bar{y}_2\mathbf{1}_{n_2} \\ \vdots \\ \bar{y}_k\mathbf{1}_{n_k} \end{bmatrix},$$

whereas the sum of squared errors is given by

$$(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2.$$

6.7 SYSTEMS OF LINEAR EQUATIONS AND THE SINGULAR VALUE DECOMPOSITION

When A is square and nonsingular, then the solution to the system of equations $A\mathbf{x} = \mathbf{c}$ can be conveniently expressed in terms of the inverse of A , as $\mathbf{x} = A^{-1}\mathbf{c}$. For this reason, it has seemed somewhat natural to deal with the solutions for the more general case in terms of the generalization of A^{-1} , A^+ . This is the approach that we have taken throughout this chapter. Alternatively, we can attack this problem by directly using the singular value decomposition, an approach that may offer more insight. In this case, we will always be able to transform our system to a simpler system of equations of the form

$$D\mathbf{y} = \mathbf{b}, \quad (6.16)$$

where \mathbf{y} is an $n \times 1$ vector of variables, \mathbf{b} is an $m \times 1$ vector of constants, and D is an $m \times n$ matrix such that $d_{ij} = 0$ if $i \neq j$. In particular, D will have one of the four forms, as given in Theorem 4.1,

$$(a) \Delta, \quad (b) [\Delta \quad (0)], \quad (c) \begin{bmatrix} \Delta \\ (0) \end{bmatrix}, \quad (d) \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix},$$

where Δ is an $r \times r$ nonsingular diagonal matrix and $r = \text{rank}(A)$. Now if D has the form given in (a), then the system (6.16) is consistent with the unique solution given by $\mathbf{y} = \Delta^{-1}\mathbf{b}$. For (b), if we partition \mathbf{y} as $\mathbf{y} = (\mathbf{y}'_1, \mathbf{y}'_2)'$, where \mathbf{y}_1 is $r \times 1$, then (6.16) reduces to

$$\Delta\mathbf{y}_1 = \mathbf{b}.$$

Thus, (6.16) is consistent and has solutions of the form

$$\mathbf{y} = \begin{bmatrix} \Delta^{-1}\mathbf{b} \\ \mathbf{y}_2 \end{bmatrix},$$

where the $(n - r) \times 1$ vector \mathbf{y}_2 is arbitrary. Since we then have $n - r$ linearly independent choices for \mathbf{y}_2 , the number of linearly independent solutions is $n - r$ if $\mathbf{b} = \mathbf{0}$ and $n - r + 1$ if $\mathbf{b} \neq \mathbf{0}$. When D has the form given in (c), the system in (6.16) takes the form

$$\begin{bmatrix} \Delta\mathbf{y} \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix},$$

where \mathbf{b}_1 is $r \times 1$ and \mathbf{b}_2 is $(m - r) \times 1$, and so it is consistent only if $\mathbf{b}_2 = \mathbf{0}$. If this is the case, the system then has a unique solution given by $\mathbf{y} = \Delta^{-1}\mathbf{b}_1$. For the final form given in (d), the system of equations in (6.16) appears as

$$\begin{bmatrix} \Delta\mathbf{y}_1 \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} \mathbf{b}_1 \\ \mathbf{b}_2 \end{bmatrix},$$

where \mathbf{y} and \mathbf{b} have been partitioned as before. As in the case of form (c), this system is consistent only if $\mathbf{b}_2 = \mathbf{0}$, and as in the case of form (b), when consistent, it has

$n - r$ linearly independent solutions if $\mathbf{b} = \mathbf{0}$ and $n - r + 1$ linearly independent solutions if $\mathbf{b} \neq \mathbf{0}$. The general solution is given by

$$\mathbf{y} = \begin{bmatrix} \Delta^{-1}\mathbf{b}_1 \\ \mathbf{y}_2 \end{bmatrix},$$

where the $(n - r) \times 1$ vector \mathbf{y}_2 is arbitrary.

All of the above can now be readily applied to the general system of equations,

$$A\mathbf{x} = \mathbf{c} \quad (6.17)$$

by utilizing the singular value decomposition of A given by $A = PDQ'$ as in Theorem 4.1. Premultiplication of this system of equations by P' produces the system of equations in (6.16), where the vector of variables is given by $\mathbf{y} = Q'\mathbf{x}$ and the vector of constants is $\mathbf{b} = P'\mathbf{c}$. Consequently, if \mathbf{y} is a solution to (6.16), then $\mathbf{x} = Q\mathbf{y}$ will be a solution to (6.17). Thus, in the case of forms (a) and (b), (6.17) is consistent with a unique solution given by

$$\mathbf{x} = Q\mathbf{y} = Q\Delta^{-1}\mathbf{b} = Q\Delta^{-1}P'\mathbf{c} = A^{-1}\mathbf{c},$$

when (a) is the form of D , whereas for form (b) the general solution is

$$\mathbf{x} = Q\mathbf{y} = [Q_1 \quad Q_2] \begin{bmatrix} \Delta^{-1}\mathbf{b} \\ \mathbf{y}_2 \end{bmatrix} = Q_1\Delta^{-1}P'\mathbf{c} + Q_2\mathbf{y}_2,$$

where Q_1 is $n \times r$ and \mathbf{y}_2 is an arbitrary $(n - r) \times 1$ vector. The term $Q_2\mathbf{y}_2$ has no effect on the value of $A\mathbf{x}$ because the columns of the $n \times (n - r)$ matrix Q_2 form a basis for the null space of A . In the case of forms (c) and (d), the system (6.17) is consistent only if $\mathbf{c} = P_1\mathbf{b}_1$, so that $P_2\mathbf{b}_2 = \mathbf{0}$, where $P = (P_1, P_2)$ and P_1 is $m \times r$; that is, because the columns of P_1 form a basis for the range of A , the system is consistent if \mathbf{c} is in the column space of A . Thus, if we partition \mathbf{c} as $\mathbf{c} = (\mathbf{c}'_1, \mathbf{c}'_2)'$, where \mathbf{c}_1 is $r \times 1$, then when form (c) holds, the unique solution will be given by

$$\mathbf{x} = Q\mathbf{y} = Q\Delta^{-1}\mathbf{b}_1 = Q\Delta^{-1}P'_1\mathbf{c}.$$

In the case of form (d), the general solution is

$$\mathbf{x} = Q\mathbf{y} = [Q_1 \quad Q_2] \begin{bmatrix} \Delta^{-1}\mathbf{b}_1 \\ \mathbf{y}_2 \end{bmatrix} = Q_1\Delta^{-1}P'_1\mathbf{c} + Q_2\mathbf{y}_2.$$

6.8 SPARSE LINEAR SYSTEMS OF EQUATIONS

The typical approach to the numerical computation of solutions to a consistent system of equations $Ax = c$, or least squares solutions when the system is inconsistent, applies some factorization of A such as the QR factorization, the singular value decomposition, or the LU decomposition, which factors A into the product of a lower triangular matrix and upper triangular matrix. Any method of this type is referred to as a direct method. One situation in which direct methods may not be appropriate is when our system of equations is large and sparse; that is, m and n are large and a relatively large number of the elements of the $m \times n$ matrix A are equal to zero. Thus, although the size of A may be quite large, its storage will not require an enormous amount of computer memory because we only need to store the nonzero values and their locations. However, when A is sparse, the factors in its decompositions need not be sparse, so if A is large enough, the computation of these factorizations may easily require more memory than is available.

If there is some particular structure to the sparsity of A , then it may be possible to implement a direct method that exploits this structure. A simple example of such a situation is one in which A is $m \times m$ and tridiagonal; that is, A has the form

$$A = \begin{bmatrix} v_1 & w_1 & 0 & \cdots & 0 & 0 & 0 \\ u_2 & v_2 & w_2 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & u_{m-1} & v_{m-1} & w_{m-1} \\ 0 & 0 & 0 & \cdots & 0 & u_m & v_m \end{bmatrix}.$$

In this case, if we define

$$L = \begin{bmatrix} r_1 & 0 & \cdots & 0 & 0 \\ u_2 & r_2 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & r_{m-1} & 0 \\ 0 & 0 & \cdots & u_m & r_m \end{bmatrix}, \quad U = \begin{bmatrix} 1 & s_1 & \cdots & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & s_{m-1} \\ 0 & 0 & \cdots & 0 & 1 \end{bmatrix},$$

where $r_1 = v_1$, $r_i = v_i - u_i w_{i-1}/r_{i-1}$, and $s_{i-1} = w_{i-1}/r_{i-1}$, for $i = 2, \dots, m$, then A can be factored as $A = LU$ as long as each $r_i \neq 0$. Thus, the two factors, L and U , are also sparse. The system $Ax = c$ can easily be solved by first solving the system $Ly = c$ and then solving the system $Ux = y$. For more details on this and adaptations of direct methods for other structured matrices, such as banded matrices and block tridiagonal matrices, see Duff, et al. (1986) and Golub and Van Loan (2013).

A second approach to the solution of sparse systems of equations uses iterative methods. In this case, a sequence of vectors, x_0, x_1, \dots is generated with x_0 being some initial vector, whereas x_j for $j = 1, 2, \dots$ is a vector that is computed using the previous vector x_{j-1} , with the property that $x_j \rightarrow x$, as $j \rightarrow \infty$, where x is the

true solution to $A\mathbf{x} = \mathbf{c}$. Typically, the computation in these methods only involves A through its product with vectors, and this is an operation that will be easy to handle if A is sparse. Two of the oldest and simplest iterative schemes are the Jacobi and Gauss–Seidel methods. If A is $m \times m$ with nonzero diagonal elements, then the system $A\mathbf{x} = \mathbf{c}$ can be written as

$$(A - D_A)\mathbf{x} + D_A\mathbf{x} = \mathbf{c},$$

which yields the identity

$$\mathbf{x} = D_A^{-1}\{\mathbf{c} - (A - D_A)\mathbf{x}\}.$$

This identity is the motivation for the Jacobi method that computes \mathbf{x}_j as

$$\mathbf{x}_j = D_A^{-1}\{\mathbf{c} - (A - D_A)\mathbf{x}_{j-1}\}.$$

On the other hand, the Gauss–Seidel method applies the splitting of A as $A = A_1 + A_2$, where A_1 is lower triangular and A_2 is upper triangular with each of its diagonal elements equal to zero. In this case, $A\mathbf{x} = \mathbf{c}$ can be rearranged as

$$A_1\mathbf{x} = \mathbf{c} - A_2\mathbf{x},$$

and this leads to the iterative scheme

$$A_1\mathbf{x}_j = \mathbf{c} - A_2\mathbf{x}_{j-1},$$

which is easily solved for \mathbf{x}_j because the system is triangular.

In recent years, some other more sophisticated iterative methods, requiring less computation and having better convergence properties, have been developed. We will briefly discuss a method for solving a system of equations which utilizes an algorithm known as the Lanczos algorithm (Lanczos, 1950). For more information on this procedure, including convergence properties, generalizations to a general $m \times n$ matrix, and to the problem of finding least squares solutions, as well as other iterative methods, the reader is referred to Young (1971), Hageman and Young (1981), and Golub and Van Loan (2013).

Consider the function

$$f(\mathbf{x}) = \frac{1}{2}\mathbf{x}'A\mathbf{x} - \mathbf{x}'\mathbf{c},$$

where \mathbf{x} is an $m \times 1$ vector and A is an $m \times m$ positive definite matrix. The vector of partial derivatives of $f(\mathbf{x})$ given by

$$\nabla f(\mathbf{x}) = \left(\frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right)' = A\mathbf{x} - \mathbf{c}$$

is sometimes referred to as the gradient of $f(\mathbf{x})$. Setting this equation equal to the zero vector, we find that the vector minimizing f , $\mathbf{x} = A^{-1}\mathbf{c}$, is also the solution to the system $A\mathbf{x} = \mathbf{c}$. Thus, a vector that approximately minimizes f will also be an approximate solution to $A\mathbf{x} = \mathbf{c}$. One iterative method for finding the minimizer \mathbf{x} involves successively finding minimizers \mathbf{x}_j of f over a j -dimensional subspace of R^m , starting with $j = 1$ and continually increasing j by 1. In particular, for some set of orthonormal $m \times 1$ vectors, $\mathbf{q}_1, \dots, \mathbf{q}_m$, we will define the j th subspace as the space with the columns of the $m \times j$ matrix, $Q_j = (\mathbf{q}_1, \dots, \mathbf{q}_j)$, as its basis. Consequently, for some $j \times 1$ vector \mathbf{y}_j ,

$$\mathbf{x}_j = Q_j \mathbf{y}_j \quad (6.18)$$

and

$$f(\mathbf{x}_j) = \min_{\mathbf{y} \in R^j} f(Q_j \mathbf{y}) = \min_{\mathbf{y} \in R^j} g(\mathbf{y}) = g(\mathbf{y}_j),$$

where

$$g(\mathbf{y}) = \frac{1}{2} \mathbf{y}' (Q_j' A Q_j) \mathbf{y} - \mathbf{y}' Q_j' \mathbf{c}.$$

Thus, the gradient of $g(\mathbf{y}_j)$ must be equal to the null vector, and so

$$(Q_j' A Q_j) \mathbf{y}_j = Q_j' \mathbf{c}. \quad (6.19)$$

To obtain \mathbf{x}_j , we can first use (6.19) to calculate \mathbf{y}_j and then use this in (6.18) to get \mathbf{x}_j . The final \mathbf{x}_j , \mathbf{x}_m , will be the solution to $A\mathbf{x} = \mathbf{c}$, but the goal here is to stop the iterative process before $j = m$ with a sufficiently accurate solution \mathbf{x}_j .

The iterative scheme described above will work with different sets of orthonormal vectors $\mathbf{q}_1, \dots, \mathbf{q}_m$, but we will see that by a judicious choice of this set, we may guarantee that the computation involved in computing the \mathbf{x}_j 's will be fairly straightforward even when A is large and sparse. These same vectors are also useful in an iterative procedure for obtaining a few of the largest and smallest eigenvalues of A . We will derive these vectors in the context of this eigenvalue problem and then later return to our discussion of the system of equations $A\mathbf{x} = \mathbf{c}$. Let λ_1 and λ_m denote the largest and smallest eigenvalues of A , whereas λ_{1j} and λ_{jj} denote the largest and smallest eigenvalues of the $j \times j$ matrix $Q_j' A Q_j$. Now we have seen in Chapter 3 that $\lambda_{1j} \leq \lambda_1$, $\lambda_{jj} \geq \lambda_m$ and that λ_1 and λ_m are the maximum and minimum values of the Rayleigh quotient,

$$R(\mathbf{x}, A) = \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}}.$$

Suppose that we have the j columns of Q_j , and we wish to find an additional vector \mathbf{q}_{j+1} so as to form the matrix Q_{j+1} and have $\lambda_{1,j+1}$ and $\lambda_{j+1,j+1}$ as close to λ_1 and λ_m as possible. If \mathbf{u}_j is a vector in the space spanned by the columns of Q_j and satisfying $R(\mathbf{u}_j, A) = \lambda_{1j}$, then because the gradient

$$\nabla R(\mathbf{u}_j, A) = \frac{2}{\mathbf{u}_j' \mathbf{u}_j} \{A \mathbf{u}_j - R(\mathbf{u}_j, A) \mathbf{u}_j\}$$

gives the direction in which $R(\mathbf{u}_j, A)$ is increasing most rapidly, we would want to choose \mathbf{q}_{j+1} so that $\nabla R(\mathbf{u}_j, A)$ is in the space spanned by the columns of Q_{j+1} . On the other hand, if \mathbf{v}_j is a vector in the space spanned by Q_j and satisfying $R(\mathbf{v}_j, A) = \lambda_{jj}$, then because $R(\mathbf{v}_j, A)$ is decreasing most rapidly in the direction given by $-\nabla R(\mathbf{v}_j, A)$, we would want to make sure that $\nabla R(\mathbf{v}_j, A)$ is also in the space spanned by the columns of Q_{j+1} . Both of these objectives can be satisfied if the columns of Q_j are spanned by the vectors $\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1$ and we select \mathbf{q}_{j+1} so that the columns of Q_{j+1} are spanned by the vectors $\mathbf{q}_1, A\mathbf{q}_1, \dots, A^j\mathbf{q}_1$ because both $\nabla R(\mathbf{u}_j, A)$ and $\nabla R(\mathbf{v}_j, A)$ are of the form $aA\mathbf{x} + b\mathbf{x}$ for some vector \mathbf{x} spanned by the columns of Q_j . Thus, we start with an initial unit vector \mathbf{q}_1 , whereas for $j \geq 2$, \mathbf{q}_j is selected as a unit vector orthogonal to $\mathbf{q}_1, \dots, \mathbf{q}_{j-1}$, and such that the columns of Q_j are spanned by the vectors $\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1$. These particular \mathbf{q}_j vectors are known as the Lanczos vectors. The calculation of the \mathbf{q}_j 's can be facilitated by the use of the tridiagonal factorization $A = PTP'$, where P is orthogonal and T has the tridiagonal form

$$T = \begin{bmatrix} \alpha_1 & \beta_1 & 0 & \cdots & 0 & 0 & 0 \\ \beta_1 & \alpha_2 & \beta_2 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & \beta_{m-2} & \alpha_{m-1} & \beta_{m-1} \\ 0 & 0 & 0 & \cdots & 0 & \beta_{m-1} & \alpha_m \end{bmatrix}.$$

Using this factorization, we find that if we choose P and \mathbf{q}_1 , so that $Pe_1 = \mathbf{q}_1$, then

$$(\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1) = P(e_1, Te_1, \dots, T^{j-1}e_1).$$

Since $(e_1, Te_1, \dots, T^{j-1}e_1)$ has upper triangular structure, the first j columns of P span the column space of $(\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1)$; that is, the \mathbf{q}_j 's can be obtained by calculating the factorization $A = PTP'$, or in other words, we can take $Q = (\mathbf{q}_1, \dots, \mathbf{q}_m) = P$. Thus, because $AQ = QT$, we have

$$A\mathbf{q}_1 = \alpha_1\mathbf{q}_1 + \beta_1\mathbf{q}_2 \quad (6.20)$$

and

$$A\mathbf{q}_j = \beta_{j-1}\mathbf{q}_{j-1} + \alpha_j\mathbf{q}_j + \beta_j\mathbf{q}_{j+1}, \quad (6.21)$$

for $j = 2, \dots, m-1$. Using these equations and the orthonormality of the \mathbf{q}_j 's, it is easily shown that $\alpha_j = \mathbf{q}_j' A \mathbf{q}_j$ for all j , and as long as $\mathbf{p}_j = (A - \alpha_j I_m)\mathbf{q}_j - \beta_{j-1}\mathbf{q}_{j-1} \neq \mathbf{0}$, then $\beta_j^2 = \mathbf{p}_j' \mathbf{p}_j$, and $\mathbf{q}_{j+1} = \mathbf{p}_j / \beta_j$ for $j = 1, \dots, m-1$, if we define $\mathbf{q}_0 = \mathbf{0}$. Thus, we can continue calculating the \mathbf{q}_j 's until we encounter a $\mathbf{p}_j = \mathbf{0}$. To see the significance of this event, let us suppose that the iterative procedure has proceeded through the first $j-1$ steps with $\mathbf{p}_i \neq \mathbf{0}$ for each $i = 2, \dots, j-1$, and so we have obtained the matrix Q_j whose columns

form a basis for the column space of $(\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1}\mathbf{q}_1)$. Note that it follows immediately from the relationship $AQ = QT$ that

$$AQ_j = Q_j T_j + \mathbf{p}_j \mathbf{e}_j',$$

where T_j is the $j \times j$ submatrix of T consisting of its first j rows and j columns. This leads to the equation $Q_j' A Q_j = T_j + Q_j' \mathbf{p}_j \mathbf{e}_j'$. However, $\mathbf{q}_i' A \mathbf{q}_i = \alpha_i$, whereas it follows from (6.20) and (6.21) that $\mathbf{q}_{i+1}' A \mathbf{q}_i = \beta_i$ and $\mathbf{q}_k' A \mathbf{q}_i = 0$ if $k > i + 1$. Thus, $Q_j' A Q_j = T_j$, and so we must have $Q_j' \mathbf{p}_j = \mathbf{0}$. Now if $\mathbf{p}_j \neq \mathbf{0}$, then $\mathbf{q}_{j+1} = \mathbf{p}_j / \beta_j$ is orthogonal to the columns of Q_j . Further, it follows from the fact that \mathbf{q}_{j+1} is a linear combination of $A\mathbf{q}_j$, \mathbf{q}_j and \mathbf{q}_{j-1} that the columns of $Q_{j+1} = (Q_j, \mathbf{q}_{j+1})$ form a basis for the column space of $(\mathbf{q}_1, A\mathbf{q}_1, \dots, A^j \mathbf{q}_1)$. If, on the other hand, $\mathbf{p}_j = \mathbf{0}$, then $AQ_j = Q_j T_j$. From this we see that the vectors $A^j \mathbf{q}_1, \dots, A^{m-1} \mathbf{q}_1$ are in the space spanned by the columns of Q_j , that is, the space spanned by the vectors $\mathbf{q}_1, A\mathbf{q}_1, \dots, A^{j-1} \mathbf{q}_1$. Consequently, the iterative procedure is complete because there are only j \mathbf{q}_i 's.

In the iterative procedure described above, the largest and smallest eigenvalues of T_j serve as approximations to the largest and smallest eigenvalues of A . In practice, the termination of this iterative process is usually not because of the encounter of a $\mathbf{p}_j = \mathbf{0}$, but because of sufficiently accurate approximations of the eigenvalues of A .

Now let us return to the problem of solving the system of equations $A\mathbf{x} = \mathbf{c}$ through the iterative procedure based on the calculation of \mathbf{y}_j in (6.19) and then \mathbf{x}_j in (6.18). We will see that the choice of the Lanczos vectors as the columns of Q_j will simplify the computations involved. For this choice of Q_j , we have already seen that $Q_j' A Q_j = T_j$, so that the system in (6.19) is a special case of the tridiagonal system of equations discussed at the beginning of this section, special in that T_j is symmetric. As a result, the matrix T_j can be factored as $T_j = L_j D_j L_j'$, where $D_j = \text{diag}(d_1, \dots, d_j)$,

$$L_j = \begin{bmatrix} 1 & 0 & \cdots & 0 & 0 \\ l_1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & \cdots & l_{j-1} & 1 \end{bmatrix},$$

$d_1 = \alpha_1$, and for $i = 2, \dots, j$, $l_{i-1} = \beta_{i-1}/d_{i-1}$ and $d_i = \alpha_i - \beta_{i-1}l_{i-1}$. Thus, the solution for \mathbf{y}_j in (6.19) can be easily found by first solving $L_j \mathbf{w}_j = Q_j' \mathbf{c}$, then $D_j \mathbf{z}_j = \mathbf{w}_j$, and finally $L_j' \mathbf{y}_j = \mathbf{z}_j$. Even as j increases, the computation required is not extensive because D_{j-1} and L_{j-1} are submatrices of D_j and L_j , and so in the j th iteration, we only need to calculate d_j and l_{j-1} to obtain D_j and L_j from D_{j-1} and L_{j-1} .

The next step is to compute \mathbf{x}_j from \mathbf{y}_j using (6.18). We will see that this also may be done with a small amount of computation. Note that if we define the $m \times j$

matrix $B_j = (\mathbf{b}_1, \dots, \mathbf{b}_j)$ so that $B_j L'_j = Q_j$, then by premultiplying the equation $T_j \mathbf{y}_j = Q'_j \mathbf{c}$ by $Q_j T_j^{-1}$ and using (6.18), we get

$$\mathbf{x}_j = Q_j T_j^{-1} Q'_j \mathbf{c} = Q_j (L_j D_j L'_j)^{-1} Q'_j \mathbf{c} = B_j \mathbf{z}_j, \quad (6.22)$$

where \mathbf{z}_j is as previously defined. It will be easier to compute \mathbf{x}_j from (6.22) than from (6.18) because B_j and \mathbf{z}_j are simple to compute after B_{j-1} and \mathbf{z}_{j-1} have already been calculated. For instance, from the definition of B_j , we see that $\mathbf{b}_1 = \mathbf{q}_1$ and $\mathbf{b}_i = \mathbf{q}_i - l_{i-1} \mathbf{b}_{i-1}$ for $i > 1$, and consequently, $B_j = (B_{j-1}, \mathbf{b}_j)$. Using the defining equations for \mathbf{w}_j and \mathbf{z}_j , we find that

$$L_j D_j \mathbf{z}_j = Q'_j \mathbf{c}. \quad (6.23)$$

If we partition \mathbf{z}_j as $\mathbf{z}_j = (\gamma'_{j-1}, \gamma_j)'$, where γ_{j-1} is a $(j-1) \times 1$ vector, then by using the fact that

$$L_j = \begin{bmatrix} L_{j-1} & \mathbf{0} \\ l_{j-1} \mathbf{e}'_{j-1} & 1 \end{bmatrix}, \quad D_j = \begin{bmatrix} D_{j-1} & \mathbf{0} \\ \mathbf{0}' & d_j \end{bmatrix},$$

we see that (6.23) implies that $L_{j-1} D_{j-1} \gamma_{j-1} = Q'_{j-1} \mathbf{c}$. As a result, $\gamma_{j-1} = \mathbf{z}_{j-1}$, and so to compute \mathbf{z}_j , we only need to compute γ_j , which is given by

$$\gamma_j = (\mathbf{q}'_j \mathbf{c} - l_{j-1} d_{j-1} \gamma_{j-1}) / d_j,$$

where γ_{j-1} is the last component of \mathbf{z}_{j-1} . Thus, (6.22) becomes

$$\begin{aligned} \mathbf{x}_j &= B_j \mathbf{z}_j = [B_{j-1} \quad \mathbf{b}_j] \begin{bmatrix} \mathbf{z}_{j-1} \\ \gamma_j \end{bmatrix} \\ &= B_{j-1} \mathbf{z}_{j-1} + \gamma_j \mathbf{b}_j = \mathbf{x}_{j-1} + \gamma_j \mathbf{b}_j, \end{aligned}$$

and so we have a simple formula for computing the j th iterative solution from \mathbf{b}_j , γ_j , and the $(j-1)$ th iterative solution \mathbf{x}_{j-1} .

PROBLEMS

6.1 Consider the system of equations $A\mathbf{x} = \mathbf{c}$, where A is the 4×3 matrix given in Problem 5.2 and

$$\mathbf{c} = \begin{bmatrix} 1 \\ 3 \\ -1 \\ 0 \end{bmatrix}.$$

(a) Show that the system is consistent.

- (b) Find a solution to this system of equations.
 (c) How many linearly independent solutions are there?

6.2 The system of equations $A\mathbf{x} = \mathbf{c}$ has A equal to the 3×4 matrix given in Problem 5.48 and

$$\mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 4 \end{bmatrix}.$$

- (a) Show that the system of equations is consistent.
 (b) Give the general solution.
 (c) Find r , the number of linearly independent solutions.
 (d) Give a set of r linearly independent solutions.

6.3 Suppose the system of equations $A\mathbf{x} = \mathbf{c}$ has

$$A = \begin{bmatrix} 5 & 2 & 1 \\ 3 & 1 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & -3 \end{bmatrix}.$$

For each \mathbf{c} given below, determine whether the system of equations is consistent.

$$(a) \quad \mathbf{c} = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad (b) \quad \mathbf{c} = \begin{bmatrix} 3 \\ 2 \\ 1 \\ -1 \end{bmatrix}, \quad (c) \quad \mathbf{c} = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}.$$

6.4 Consider the system of equations $A\mathbf{x} = \mathbf{c}$, where

$$A = \begin{bmatrix} 1 & 1 & -1 & 0 & 2 \\ 2 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 3 \\ 1 \end{bmatrix}.$$

- (a) Show that the system of equations is consistent.
 (b) Give the general solution.
 (c) Find r , the number of linearly independent solutions.
 (d) Give a set of r linearly independent solutions.

6.5 Prove Theorem 6.5.

6.6 Consider the system of equations $AXB = C$, where X is a 3×3 matrix of variables and

$$A = \begin{bmatrix} 1 & 3 & 1 \\ 3 & 2 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & -1 \\ 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 4 & 2 \\ 2 & 1 \end{bmatrix}.$$

- (a) Show that the system of equations is consistent.
 (b) Find the form of the general solution to this system.
- 6.7** The general solution of a consistent system of equations was given in Theorem 6.4 as $A^-c + (I_n - A^-A)y$. Show that the two vectors A^-c and $(I_n - A^-A)y$ are linearly independent if $c \neq 0$ and $(I_n - A^-A)y \neq 0$.
- 6.8** Suppose we want to find solutions x to (6.1) under the restriction that the solutions are in some vector space S . Let P_S be the projection matrix for the orthogonal projection onto S . Show that a restricted solution exists if $AP_S(AP_S)^-c = c$ and, in this case, the general solution is given by

$$x_y = P_S(AP_S)^-c + P_S\{I_n - (AP_S)^-AP_S\}y,$$

where y is an arbitrary $n \times 1$ vector.

- 6.9** Suppose the $m \times n$ matrix A and $m \times 1$ vector $c \neq 0$ are such that A^-c is the same for all choices of A^- . Use Theorem 5.24 to show that if $Ax = c$ is a consistent system of equations, then it has a unique solution.
- 6.10** Let A be an $m \times n$ matrix, whereas c and d are $m \times 1$ and $n \times 1$ vectors, respectively.
- (a) Show that $BAX = BC$ has the same set of solutions for the $n \times 1$ vector x as $Ax = c$ if the $p \times m$ matrix B has full column rank.
- (b) Show that $A'AX = A'Ad$ has the same set of solutions for x as $Ax = Ad$.
- 6.11** For the homogeneous system of equations $Ax = 0$ in which

$$A = \begin{bmatrix} -1 & 3 & -2 & 1 \\ 2 & -3 & 0 & -2 \end{bmatrix},$$

determine r , the number of linearly independent solutions, and find a set of r linearly independent solutions.

- 6.12** Suppose that x_* is a solution to the system of equations $Ax = c$ and y_* is a solution to the system of equations $A'y = d$. Show that $d'x_* = c'y_*$.
- 6.13** Suppose $Ax = c$ is a consistent system of equations, and let G be a matrix that satisfies conditions 1 and 4 of the Moore–Penrose generalized inverse of A . Show that Gc is the minimal solution of the system $Ax = c$. That is, show that for any other solution $x_* \neq Gc$, $x'_*x_* > c'G'Gc$.
- 6.14** Show that if the system of equations $AXB = C$ is consistent, then the solution is unique if and only if A has full column rank and B has full row rank.
- 6.15** Suppose A is an $m \times m$ matrix satisfying $AX_1 = X_1\Lambda_1$ for some $m \times r$ full rank matrix X_1 and diagonal matrix $\Lambda_1 = \text{diag}(\lambda_1, \dots, \lambda_r)$. If X_2 is an $m \times (m - r)$ matrix such that $X = [X_1 \ X_2]$ is nonsingular, show that A can be expressed as $A = X_1\Lambda_1X_1^- - WX_2^-$, where W is an $m \times (m - r)$ matrix, and X_1^- and X_2^- are the generalized inverses of X_1 and X_2 satisfying $X^{-1'} = [X_1^{-1'} \ X_2^{-1'}]$.

6.16 Let

$$A = \begin{bmatrix} 1 & -1 & 1 & 1 \\ 2 & 3 & 1 & -1 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 1 \\ 2 \end{bmatrix},$$

$$B = \begin{bmatrix} 2 & 1 & 2 & -1 \\ 0 & 1 & 1 & 1 \end{bmatrix}, \quad \mathbf{d} = \begin{bmatrix} 2 \\ 4 \end{bmatrix}.$$

- (a) Show that the system $A\mathbf{x} = \mathbf{c}$ is consistent and has three linearly independent solutions.
- (b) Show that the system $B\mathbf{x} = \mathbf{d}$ is consistent and has three linearly independent solutions.
- (c) Show that the systems $A\mathbf{x} = \mathbf{c}$ and $B\mathbf{x} = \mathbf{d}$ have a common solution and that this common solution is unique.
- 6.17** Consider the systems of equations $AX = C$ and $XB = D$, where A is $m \times n$, B is $p \times q$, C is $m \times p$, and D is $n \times q$.
- (a) Show that the two systems of equations have a common solution X if and only if each system is consistent and $AD = CB$.
- (b) Show that the general common solution is given by

$$X_* = A^-C + (I_n - A^-A)DB^- + (I_n - A^-A)Y(I_p - BB^-),$$

where Y is an arbitrary $n \times p$ matrix.

- 6.18** Suppose the $r \times m$ matrix A_1 is a generalized inverse of the $m \times r$ matrix X_1 , which has rank r . We wish to find $m \times (m - r)$ matrices X_2 and A'_2 such that $X = [X_1 \ X_2]$ is nonsingular and $X^{-1} = [A'_1 \ A'_2]'$. Show that X_2 and A_2 can be expressed as $X_2 = (I_m - X_1A_1)Y$ and $A_2 = \{(I_m - X_1A_1)Y\}^-(I_m - X_1A_1)$, where Y is an arbitrary $m \times (m - r)$ matrix for which $\text{rank}\{(I_m - X_1A_1)Y\} = m - r$.
- 6.19** In Problem 5.49, a least squares inverse was found for the matrix

$$A = \begin{bmatrix} 1 & -1 & -2 & 1 \\ -2 & 4 & 3 & -2 \\ 1 & 1 & -3 & 1 \end{bmatrix}.$$

- (a) Use this least squares inverse to show that the system of equations $A\mathbf{x} = \mathbf{c}$ is inconsistent, where $\mathbf{c}' = (2, 1, 5)$.
- (b) Find a least squares solution.
- (c) Compute the sum of squared errors for a least squares solution to this system of equations.

6.20 Consider the system of equations $A\mathbf{x} = \mathbf{c}$, where

$$A = \begin{bmatrix} 1 & 0 & 2 \\ 2 & -1 & 3 \\ -1 & 2 & 0 \\ -2 & 1 & -3 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 2 \\ 2 \\ 5 \\ 0 \end{bmatrix}.$$

- (a) Find a least squares inverse of A .
 - (b) Show that the system of equations is inconsistent.
 - (c) Find a least squares solution.
 - (d) Is this solution unique?
- 6.21** Show that \mathbf{x}_* is a least squares solution to the system of equations $A\mathbf{x} = \mathbf{c}$ if and only if

$$A' A \mathbf{x}_* = A' \mathbf{c}.$$

6.22 Let A be an $m \times n$ matrix, and \mathbf{x}_* , \mathbf{y}_* , and \mathbf{c} be $n \times 1$, $m \times 1$, and $m \times 1$ vectors, respectively. Suppose that \mathbf{x}_* and \mathbf{y}_* are such that the system of equations

$$\begin{bmatrix} I_m & A \\ A' & (0) \end{bmatrix} \begin{bmatrix} \mathbf{y}_* \\ \mathbf{x}_* \end{bmatrix} = \begin{bmatrix} \mathbf{c} \\ \mathbf{0} \end{bmatrix}$$

holds. Show that \mathbf{x}_* then must be a least squares solution to the system $A\mathbf{x} = \mathbf{c}$.

6.23 The balanced two-way classification model with interaction is of the form

$$y_{ijk} = \mu + \tau_i + \gamma_j + \eta_{ij} + \epsilon_{ijk},$$

where $i = 1, \dots, a$, $j = 1, \dots, b$, and $k = 1, \dots, n$. The parameter μ represents an overall effect, τ_i is an effect due to the i th level of factor one, γ_j is an effect due to the j th level of factor two, and η_{ij} is an effect due to the interaction of the i th and j th levels of factors one and two; as usual, the ϵ_{ijk} 's represent independent random errors, each distributed as $N(0, \sigma^2)$.

- (a) Set up the vectors \mathbf{y} , $\boldsymbol{\beta}$, and $\boldsymbol{\epsilon}$ and the matrix X so that the two-way model above can be written in the matrix form $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$.
 - (b) Find the rank, r , of X . Determine a set of r linearly independent estimable functions of the parameters, μ , τ_i , γ_j , and η_{ij} .
 - (c) Find a least squares solution for the parameter vector $\boldsymbol{\beta}$.
- 6.24** Consider the regression model

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where X is $N \times m$, $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 C)$, and C is a known positive definite matrix. In Example 4.6, for the case in which X is full column rank, we

obtained the generalized least squares estimator, $\hat{\beta} = (X'C^{-1}X)^{-1}X'C^{-1}\mathbf{y}$, that minimizes

$$(\mathbf{y} - X\hat{\beta})'C^{-1}(\mathbf{y} - X\hat{\beta}). \quad (6.24)$$

Show that if X is less than full column rank, then the generalized least squares solution for β that minimizes (6.24) is given by

$$\hat{\beta} = (X'C^{-1}X)^{-}X'C^{-1}\mathbf{y} + \{I_m - (X'C^{-1}X)^{-}X'C^{-1}X\}\mathbf{u},$$

where \mathbf{u} is an arbitrary $m \times 1$ vector.

6.25 Restricted least squares obtains the vectors $\hat{\beta}$ that minimize

$$(\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}),$$

subject to the restriction that $\hat{\beta}$ satisfies $B\hat{\beta} = \mathbf{b}$, where B is $p \times m$ and \mathbf{b} is $p \times 1$, such that $BB^{-}\mathbf{b} = \mathbf{b}$. Use Theorem 6.4 to find the general solution $\hat{\beta}_{\mathbf{u}}$ to the consistent system of equations $B\hat{\beta} = \mathbf{b}$, where $\hat{\beta}_{\mathbf{u}}$ depends on an arbitrary vector \mathbf{u} . Substitute this expression for $\hat{\beta}$ into $(\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta})$, and then use Theorem 6.14 to obtain the general least squares solution $\hat{\beta}_{\omega}$, for \mathbf{u} , where \mathbf{u}_{ω} depends on an arbitrary vector \mathbf{w} . Substitute \mathbf{u}_{ω} for \mathbf{u} in $\hat{\beta}_{\mathbf{u}}$ to show that the general restricted least squares solution for β is given by

$$\begin{aligned} \hat{\beta}_{\omega} &= B^{-}\mathbf{b} + (I_m - B^{-}B)\{[X(I_m - B^{-}B)]^L(\mathbf{y} - XB^{-}\mathbf{b}) \\ &\quad + (I_m - [X(I_m - B^{-}B)]^LX(I_m - B^{-}B))\mathbf{w}\}. \end{aligned}$$

6.26 In the previous problem, show that if we use the Moore–Penrose inverse as the least squares inverse of $[X(I_m - B^{-}B)]$ in the expression given for $\hat{\beta}_{\omega}$, then it simplifies to

$$\begin{aligned} \hat{\beta}_{\omega} &= B^{-}\mathbf{b} + [X(I_m - B^{-}B)]^{+}(\mathbf{y} - XB^{-}\mathbf{b}) \\ &\quad + (I_m - B^{-}B)\{I_m - [X(I_m - B^{-}B)]^{+}X(I_m - B^{-}B)\}\mathbf{w}. \end{aligned}$$

7

PARTITIONED MATRICES

7.1 INTRODUCTION

The concept of partitioning matrices was first introduced in Chapter 1, and we have subsequently used partitioned matrices throughout this text. Up to this point, most of our applications involving partitioned matrices have utilized only the simple operations of matrix addition and matrix multiplication. In this chapter, we will obtain expressions for such things as the inverse, determinant, and rank of a matrix in terms of its submatrices. We will restrict attention to an $m \times m$ matrix A that is partitioned into the 2×2 form

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad (7.1)$$

where A_{11} is $m_1 \times m_1$, A_{12} is $m_1 \times m_2$, A_{21} is $m_2 \times m_1$, and A_{22} is $m_2 \times m_2$. Additional results, including ones for which A_{11} and A_{22} are not square matrices, can be found in Harville (1997).

7.2 THE INVERSE

In this section, we obtain an expression for the inverse of A when it and at least one of the submatrices down the diagonal are nonsingular.

Theorem 7.1 Let the $m \times m$ matrix A be partitioned as in (7.1), and suppose that A is nonsingular. For notational convenience, write $B = A^{-1}$ and partition B as

$$B = \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix},$$

where the submatrices of B are of the same sizes as the corresponding submatrices of A . Then if A_{11} and $A_{22} - A_{21}A_{11}^{-1}A_{12}$ are nonsingular, we have

- (a) $B_{11} = A_{11}^{-1} + A_{11}^{-1}A_{12}B_{22}A_{21}A_{11}^{-1}$,
- (b) $B_{22} = (A_{22} - A_{21}A_{11}^{-1}A_{12})^{-1}$,
- (c) $B_{12} = -A_{11}^{-1}A_{12}B_{22}$,
- (d) $B_{21} = -B_{22}A_{21}A_{11}^{-1}$,

whereas if A_{22} and $A_{11} - A_{12}A_{22}^{-1}A_{21}$ are nonsingular, we have

- (e) $B_{11} = (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1}$,
- (f) $B_{22} = A_{22}^{-1} + A_{22}^{-1}A_{21}B_{11}A_{12}A_{22}^{-1}$,
- (g) $B_{12} = -B_{11}A_{12}A_{22}^{-1}$,
- (h) $B_{21} = -A_{22}^{-1}A_{21}B_{11}$.

Proof. Suppose that A_{11} and $A_{22} - A_{21}A_{11}^{-1}A_{12}$ are nonsingular. The matrix equation

$$AB = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} B_{11} & B_{12} \\ B_{21} & B_{22} \end{bmatrix} = \begin{bmatrix} I_{m_1} & (0) \\ (0) & I_{m_2} \end{bmatrix} = I_m$$

yields the four equations

$$A_{11}B_{11} + A_{12}B_{21} = I_{m_1}, \quad (7.2)$$

$$A_{21}B_{12} + A_{22}B_{22} = I_{m_2}, \quad (7.3)$$

$$A_{11}B_{12} + A_{12}B_{22} = (0), \quad (7.4)$$

$$A_{21}B_{11} + A_{22}B_{21} = (0). \quad (7.5)$$

Solving (7.4) for B_{12} immediately leads to the expression given in (c) for B_{12} . Substituting this solution for B_{12} into (7.3) and solving for B_{22} yields the expression given for B_{22} in (b). From (7.2) we get

$$B_{11} = A_{11}^{-1} - A_{11}^{-1}A_{12}B_{21}. \quad (7.6)$$

Substituting this result in (7.5) yields

$$A_{21}A_{11}^{-1} - A_{21}A_{11}^{-1}A_{12}B_{21} + A_{22}B_{21} = (0),$$

from which we get the expression for B_{21} given in (d). Finally, when this result is substituted back in (7.6), we obtain the expression given for B_{11} in (a). The expressions given in (e)–(h) are obtained in a similar fashion. \square

Note that, in general, A_{11} and A_{22} do not need to be nonsingular for A to be nonsingular. For instance, if $m_1 = m_2$, $A_{11} = A_{22} = (0)$, and A_{12} and A_{21} are nonsingular, it is easily verified that

$$A^{-1} = \begin{bmatrix} (0) & A_{21}^{-1} \\ A_{12}^{-1} & (0) \end{bmatrix}.$$

Theorem 7.1 and additional results to be developed later in this chapter illustrate the importance of the matrices $A_{22} - A_{21}A_{11}^{-1}A_{12}$ and $A_{11} - A_{12}A_{22}^{-1}A_{21}$ in the analysis of the partitioned matrix A given in (7.1). These matrices are commonly referred to as Schur complements. In particular, $A_{22} - A_{21}A_{11}^{-1}A_{12}$ is called the Schur complement of A_{11} in A , whereas $A_{11} - A_{12}A_{22}^{-1}A_{21}$ is called the Schur complement of A_{22} in A . Some results involving Schur complements in addition to the ones given in this chapter can be found in Ouellette (1981) and Zhang (2005).

Example 7.1 Consider the regression model

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where \mathbf{y} is $N \times 1$, X is $N \times (k+1)$, $\boldsymbol{\beta}$ is $(k+1) \times 1$, and $\boldsymbol{\epsilon}$ is $N \times 1$. Suppose that $\boldsymbol{\beta}$ and X are partitioned as $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1, \boldsymbol{\beta}'_2)'$ and $X = (X_1, X_2)$ so that the product $X_1\boldsymbol{\beta}_1$ is defined, and we are interested in comparing this complete regression model with the reduced regression model,

$$\mathbf{y} = X_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}.$$

If X has full column rank, then the least squares estimators for the two models are $\hat{\boldsymbol{\beta}} = (X'X)^{-1}X'\mathbf{y}$ and $\hat{\boldsymbol{\beta}}_1 = (X'_1X_1)^{-1}X'_1\mathbf{y}$, respectively, and the difference in the sums of squared errors for the two models,

$$\begin{aligned} & (\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1)'(\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1) - (\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}) \\ &= \mathbf{y}'(I_N - X_1(X'_1X_1)^{-1}X'_1)\mathbf{y} - \mathbf{y}'(I_N - X(X'X)^{-1}X')\mathbf{y} \\ &= \mathbf{y}'X(X'X)^{-1}X'\mathbf{y} - \mathbf{y}'X_1(X'_1X_1)^{-1}X'_1\mathbf{y}, \end{aligned} \tag{7.7}$$

gives the reduction in the sum of squared errors attributable to the inclusion of the term $X_2\boldsymbol{\beta}_2$ in the complete model. By using the geometrical properties of least squares regression in Example 2.11, we showed that this reduction in the sum of squared errors simplifies to

$$\mathbf{y}'X_{2*}(X'_{2*}X_{2*})^{-1}X'_{2*}\mathbf{y},$$

where $X_{2*} = (I_N - X_1(X_1'X_1)^{-1}X_1')X_2$. An alternative way of showing this, which we illustrate here, uses Theorem 7.1. Now $X'X$ can be partitioned as

$$X'X = \begin{bmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{bmatrix},$$

and so if we let

$$C = (X_2'X_2 - X_2'X_1(X_1'X_1)^{-1}X_1'X_2)^{-1} = (X_{2*}'X_{2*})^{-1},$$

we find from a direct application of Theorem 7.1 that

$$(X'X)^{-1} = \begin{bmatrix} (X_1'X_1)^{-1} + (X_1'X_1)^{-1}X_1'X_2CX_2'X_1(X_1'X_1)^{-1} & -(X_1'X_1)^{-1}X_1'X_2C \\ -CX_2'X_1(X_1'X_1)^{-1} & C \end{bmatrix}.$$

Substituting this into (7.7) and then simplifying, we get $\mathbf{y}'X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}'\mathbf{y}$, as required.

7.3 THE DETERMINANT

In this section, we will begin by obtaining an expression for the determinant of A when at least one of the submatrices A_{11} and A_{22} is nonsingular. Before doing this, we will first consider some special cases.

Theorem 7.2 Let the $m \times m$ matrix A be partitioned as in (7.1). If $A_{22} = I_{m_2}$, and $A_{12} = (0)$ or $A_{21} = (0)$, then $|A| = |A_{11}|$.

Proof. To find the determinant

$$|A| = \begin{vmatrix} A_{11} & (0) \\ A_{21} & I_{m_2} \end{vmatrix},$$

first apply the cofactor expansion formula for a determinant on the last column of A to obtain

$$|A| = \begin{vmatrix} A_{11} & (0) \\ B & I_{m_2-1} \end{vmatrix},$$

where B is the $(m_2 - 1) \times m_1$ matrix obtained by deleting the last row from A_{21} . Repeating this process another $(m_2 - 1)$ times yields $|A| = |A_{11}|$. In a similar fashion, we obtain $|A| = |A_{11}|$ when $A_{21} = (0)$, by repeatedly expanding along the last row. \square

Clearly we have a result analogous to Theorem 7.2 when $A_{11} = I_{m_1}$ and $A_{12} = (0)$ or $A_{21} = (0)$. Also, Theorem 7.2 can be generalized to Theorem 7.3.

Theorem 7.3 Let the $m \times m$ matrix A be partitioned as in (7.1). If $A_{12} = (0)$ or $A_{21} = (0)$, then $|A| = |A_{11}||A_{22}|$.

Proof. Observe that

$$|A| = \begin{vmatrix} A_{11} & (0) \\ A_{21} & A_{22} \end{vmatrix} = \begin{vmatrix} A_{11} & (0) \\ A_{21} & I_{m_2} \end{vmatrix} \begin{vmatrix} I_{m_1} & (0) \\ (0) & A_{22} \end{vmatrix} = |A_{11}||A_{22}|,$$

where the last equality follows from Theorem 7.2. A similar proof yields $|A| = |A_{11}||A_{22}|$ when $A_{21} = (0)$. \square

We are now ready to find an expression for the determinant of A when the only thing we know is that A_{11} or A_{22} is nonsingular.

Theorem 7.4 Let the $m \times m$ matrix A be partitioned as in (7.1). Then

- (a) $|A| = |A_{22}||A_{11} - A_{12}A_{22}^{-1}A_{21}|$, if A_{22} is nonsingular,
- (b) $|A| = |A_{11}||A_{22} - A_{21}A_{11}^{-1}A_{12}|$, if A_{11} is nonsingular.

Proof. Suppose that A_{22} is nonsingular. Note that, in this case, the identity

$$\begin{aligned} \begin{bmatrix} I_{m_1} & -A_{12}A_{22}^{-1} \\ (0) & I_{m_2} \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} I_{m_1} & (0) \\ -A_{22}^{-1}A_{21} & I_{m_2} \end{bmatrix} \\ = \begin{bmatrix} A_{11} - A_{12}A_{22}^{-1}A_{21} & (0) \\ (0) & A_{22} \end{bmatrix} \end{aligned}$$

holds. After taking the determinant of both sides of this identity and using the previous theorem, we immediately get (a). The proof of (b) is obtained in a similar fashion by using the identity

$$\begin{aligned} \begin{bmatrix} I_{m_1} & (0) \\ -A_{21}A_{11}^{-1} & I_{m_2} \end{bmatrix} \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \begin{bmatrix} I_{m_1} - A_{11}^{-1}A_{12} \\ (0) & I_{m_2} \end{bmatrix} \\ = \begin{bmatrix} A_{11} & (0) \\ (0) & A_{22} - A_{21}A_{11}^{-1}A_{12} \end{bmatrix}. \end{aligned}$$

\square

Example 7.2 We will find the determinant and inverse of the $2m \times 2m$ matrix A given by

$$A = \begin{bmatrix} aI_m & \mathbf{1}_m \mathbf{1}_m' \\ \mathbf{1}_m \mathbf{1}_m' & bI_m \end{bmatrix},$$

where a and b are nonzero scalars. Using (a) of Theorem 7.4, we find that

$$\begin{aligned} |A| &= |bI_m| |aI_m - \mathbf{1}_m \mathbf{1}_m' (bI_m)^{-1} \mathbf{1}_m \mathbf{1}_m'| \\ &= b^m \left| aI_m - \frac{m}{b} \mathbf{1}_m \mathbf{1}_m' \right| \\ &= b^m a^{m-1} \left(a - \frac{m^2}{b} \right), \end{aligned}$$

where we have used the result of Problem 3.29(e) in the last step. The matrix A will be nonsingular if $|A| \neq 0$ or, equivalently, if

$$a \neq \frac{m^2}{b}.$$

In this case, using Theorem 7.1, we find that

$$\begin{aligned} B_{11} &= (aI_m - \mathbf{1}_m \mathbf{1}_m' (bI_m)^{-1} \mathbf{1}_m \mathbf{1}_m')^{-1} \\ &= \left(aI_m - \frac{m}{b} \mathbf{1}_m \mathbf{1}_m' \right)^{-1} \\ &= a^{-1} I_m + \left\{ \frac{m}{a(ab - m^2)} \right\} \mathbf{1}_m \mathbf{1}_m', \end{aligned}$$

where this last expression follows from Problem 3.29(d). In a similar fashion, we find that

$$\begin{aligned} B_{22} &= (bI_m - \mathbf{1}_m \mathbf{1}_m' (aI_m)^{-1} \mathbf{1}_m \mathbf{1}_m')^{-1} \\ &= \left(bI_m - \frac{m}{a} \mathbf{1}_m \mathbf{1}_m' \right)^{-1} \\ &= b^{-1} I_m + \left\{ \frac{m}{b(ab - m^2)} \right\} \mathbf{1}_m \mathbf{1}_m'. \end{aligned}$$

The remaining submatrices of $B = A^{-1}$ are given by

$$\begin{aligned} B_{12} &= -(aI_m)^{-1} \mathbf{1}_m \mathbf{1}_m' \left(b^{-1} I_m + \left\{ \frac{m}{b(ab - m^2)} \right\} \mathbf{1}_m \mathbf{1}_m' \right) \\ &= -(ab - m^2)^{-1} \mathbf{1}_m \mathbf{1}_m', \end{aligned}$$

and because A is symmetric, $B_{21} = B'_{12} = B_{12}$. Putting this all together, we have

$$A^{-1} = B = \begin{bmatrix} a^{-1}(I_m + mc\mathbf{1}_m \mathbf{1}_m') & -c\mathbf{1}_m \mathbf{1}_m' \\ -c\mathbf{1}_m \mathbf{1}_m' & b^{-1}(I_m + mc\mathbf{1}_m \mathbf{1}_m') \end{bmatrix},$$

where $c = (ab - m^2)^{-1}$.

Theorem 7.4 can be used to prove the following useful inequality.

Theorem 7.5 Suppose the $m \times m$ matrix A is partitioned as in (7.1) and is a positive definite matrix. Then

$$|A| \leq |A_{11}| |A_{22}|.$$

Proof. Since $A_{21} A_{11}^{-1} A_{12} = A'_{12} A_{11}^{-1} A_{12}$ is nonnegative definite, it follows from Theorem 3.28 that

$$\begin{aligned} \lambda_h(A_{22}) &= \lambda_h(A_{22} - A'_{12} A_{11}^{-1} A_{12} + A'_{12} A_{11}^{-1} A_{12}) \\ &\geq \lambda_h(A_{22} - A'_{12} A_{11}^{-1} A_{12}) \end{aligned}$$

for $h = 1, \dots, m_2$. Consequently, $|A_{22}| \geq |A_{22} - A'_{12} A_{11}^{-1} A_{12}|$, and so the result follows from Theorem 7.4(b). \square

Theorem 7.1 and Theorem 7.4 are helpful in showing that conditional distributions formed from a random vector having a multivariate normal distribution are also multivariate normal.

Example 7.3 Suppose that $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where the covariance matrix Ω is positive definite. Partition \mathbf{x} as $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)'$, where \mathbf{x}_1 is $m_1 \times 1$ and \mathbf{x}_2 is $m_2 \times 1$. Let $\boldsymbol{\mu} = (\boldsymbol{\mu}'_1, \boldsymbol{\mu}'_2)'$ be partitioned similar to \mathbf{x} , whereas

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega'_{12} & \Omega_{22} \end{bmatrix},$$

where Ω_{11} is $m_1 \times m_1$ and Ω_{22} is $m_2 \times m_2$. We wish to determine the conditional distribution of \mathbf{x}_1 given \mathbf{x}_2 . We will do this by obtaining the conditional density of \mathbf{x}_1 given \mathbf{x}_2 , which is defined by

$$f_{1|2}(\mathbf{x}_1 | \mathbf{x}_2) = \frac{f(\mathbf{x})}{f_2(\mathbf{x}_2)},$$

where $f(\mathbf{x})$ and $f_2(\mathbf{x}_2)$ are the density functions of \mathbf{x} and \mathbf{x}_2 , respectively. Since $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, which implies that $\mathbf{x}_2 \sim N_{m_2}(\boldsymbol{\mu}_2, \Omega_{22})$, it immediately follows that

$$f(\mathbf{x}) = \frac{1}{(2\pi)^{m/2} |\Omega|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \Omega^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

and

$$f_2(\mathbf{x}_2) = \frac{1}{(2\pi)^{m_2/2} |\Omega_{22}|^{1/2}} \exp \left\{ -\frac{1}{2} (\mathbf{x}_2 - \boldsymbol{\mu}_2)' \Omega_{22}^{-1} (\mathbf{x}_2 - \boldsymbol{\mu}_2) \right\}.$$

As a result, $f_{1|2}(\mathbf{x}_1 | \mathbf{x}_2)$ has the form

$$f_{1|2}(\mathbf{x}_1 | \mathbf{x}_2) = \frac{1}{(2\pi)^{m_1/2} a} e^{-b/2}.$$

Now using Theorem 7.4(a), we find that

$$\begin{aligned} a &= |\Omega|^{1/2} |\Omega_{22}|^{-1/2} \\ &= |\Omega_{22}|^{1/2} |B_1|^{1/2} |\Omega_{22}|^{-1/2} = |B_1|^{1/2}, \end{aligned}$$

where B_1 is the Schur complement $\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega'_{12}$. Also from Theorem 7.1(e)–(h),

$$\Omega^{-1} = \begin{bmatrix} B_1^{-1} & -B_1^{-1}\Omega_{12}\Omega_{22}^{-1} \\ -\Omega_{22}^{-1}\Omega'_{12}B_1^{-1} & \Omega_{22}^{-1} + \Omega_{22}^{-1}\Omega'_{12}B_1^{-1}\Omega_{12}\Omega_{22}^{-1} \end{bmatrix},$$

so we have

$$\begin{aligned} b &= (\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu}) - (\mathbf{x}_2 - \boldsymbol{\mu}_2)'\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &= (\mathbf{x}_1 - \boldsymbol{\mu}_1)'B_1^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1) + (\mathbf{x}_2 - \boldsymbol{\mu}_2)'\Omega_{22}^{-1}\Omega'_{12}B_1^{-1}\Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &\quad - (\mathbf{x}_1 - \boldsymbol{\mu}_1)'B_1^{-1}\Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2) \\ &\quad - (\mathbf{x}_2 - \boldsymbol{\mu}_2)'\Omega_{22}^{-1}\Omega'_{12}B_1^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1) \\ &= \{\mathbf{x}_1 - \boldsymbol{\mu}_1 - \Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)\}'B_1^{-1}\{\mathbf{x}_1 - \boldsymbol{\mu}_1 - \Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)\} \\ &= (\mathbf{x}_1 - \mathbf{c})'B_1^{-1}(\mathbf{x}_1 - \mathbf{c}), \end{aligned}$$

where $\mathbf{c} = \boldsymbol{\mu}_1 + \Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2)$. Thus,

$$f_{1|2}(\mathbf{x}_1|\mathbf{x}_2) = \frac{1}{(2\pi)^{m_1/2}|B_1|^{1/2}} \exp \left\{ -\frac{1}{2}(\mathbf{x}_1 - \mathbf{c})'B_1^{-1}(\mathbf{x}_1 - \mathbf{c}) \right\},$$

which is the $N_{m_1}(\mathbf{c}, B_1)$ density function; that is, we have shown that

$$\mathbf{x}_1|\mathbf{x}_2 \sim N_{m_1}(\boldsymbol{\mu}_1 + \Omega_{12}\Omega_{22}^{-1}(\mathbf{x}_2 - \boldsymbol{\mu}_2), \Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega'_{12}).$$

Theorem 7.4 can be used to prove the following theorem which gives necessary and sufficient conditions for an $m \times m$ symmetric matrix to be positive definite.

Theorem 7.6 Let A be an $m \times m$ symmetric matrix, and let A_k be its leading $k \times k$ principal submatrix; that is, A_k is the matrix obtained by deleting the last $m - k$ rows and columns of A . Then A is positive definite if and only if all of its leading principal minors, $|A_1|, \dots, |A_m|$, are positive.

Proof. Suppose first that A is positive definite so that $|A_m| > 0$. For an arbitrary $k \times 1$ vector $\mathbf{x} \neq \mathbf{0}$ where $k < m$, define the $m \times 1$ vector $\mathbf{y} = (\mathbf{x}', \mathbf{0}')'$. Then because $\mathbf{y} \neq \mathbf{0}$ and A is positive definite, we must have

$$\mathbf{y}'A\mathbf{y} = \mathbf{x}'A_k\mathbf{x} > 0.$$

Thus, A_k is positive definite, and so $|A_k| > 0$. Next assume that $|A_k| > 0$ for $k = 1, \dots, m$. We will use induction to prove that $A = A_m$ is positive definite. Trivially, A_1 is positive definite because $|A_1| > 0$ and A_1 is 1×1 . If A_k is positive definite for some $k < m$, we will show that A_{k+1} must also be positive definite. For some $k \times 1$ vector \mathbf{b} and scalar c , A_{k+1} can be partitioned into the form

$$A_{k+1} = \begin{bmatrix} A_k & \mathbf{b} \\ \mathbf{b}' & c \end{bmatrix}.$$

Note that

$$B' A_{k+1} B = \begin{bmatrix} A_k & \mathbf{0} \\ \mathbf{0}' & c - \mathbf{b}' A_k^{-1} \mathbf{b} \end{bmatrix},$$

where

$$B = \begin{bmatrix} I_k & -A_k^{-1} \mathbf{b} \\ \mathbf{0}' & 1 \end{bmatrix},$$

and A_{k+1} is positive definite if and only if $B' A_{k+1} B$ is positive definite. Since A_k is positive definite, A_{k+1} is positive definite if and only if $c - \mathbf{b}' A_k^{-1} \mathbf{b}$ is positive. However, we know that $|A_k| > 0$, $|A_{k+1}| > 0$, and from Theorem 7.4

$$|A_{k+1}| = |A_k|(c - \mathbf{b}' A_k^{-1} \mathbf{b}),$$

so the result follows. □

We will also use Theorem 7.4 to establish Theorem 7.7.

Theorem 7.7 Let A and B be $m \times n$ and $n \times m$ matrices, respectively. Then

$$|I_m + AB| = |I_n + BA|.$$

Proof. Note that

$$\begin{bmatrix} I_m & A \\ -B & I_n \end{bmatrix} \begin{bmatrix} I_m & (0) \\ B & I_n \end{bmatrix} = \begin{bmatrix} I_m + AB & A \\ (0) & I_n \end{bmatrix},$$

so that by taking the determinant of both sides and using Theorem 7.4, we obtain the identity

$$\begin{vmatrix} I_m & A \\ -B & I_n \end{vmatrix} = |I_m + AB|. \quad (7.8)$$

Similarly, observe that

$$\begin{bmatrix} I_m & (0) \\ B & I_n \end{bmatrix} \begin{bmatrix} I_m & A \\ -B & I_n \end{bmatrix} = \begin{bmatrix} I_m & A \\ (0) & I_n + BA \end{bmatrix},$$

so that

$$\begin{vmatrix} I_m & A \\ -B & I_n \end{vmatrix} = |I_n + BA|. \quad (7.9)$$

The result now follows by equating (7.8) and (7.9). \square

Corollary 7.7.1 follows directly from Theorem 7.7 if we replace A by $-\lambda^{-1}A$.

Corollary 7.7.1 Let A and B be $m \times n$ and $n \times m$ matrices. Then the nonzero eigenvalues of AB are the same as the nonzero eigenvalues of BA .

Suppose that both A_{11} and A_{22} are singular so that we cannot use Theorem 7.4 to compute the determinant of A . A natural question in this situation is whether the formulas in Theorem 7.4 still hold if we replace the inverses by generalized inverses. A simple example will illustrate that, in general, this is not the case. For instance, if

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix},$$

then $|A| = -1$. However, the formulas $|a_{11}||a_{22} - a_{21}a_{11}^-a_{12}|$ and $|a_{22}||a_{11} - a_{12}a_{22}^-a_{21}|$ both yield 0 regardless of the choice of the generalized inverses a_{11}^- and a_{22}^- .

Conditions under which we can replace the inverses in Theorem 7.4 by generalized inverses are given in our next theorem. The matrices $A_{22} - A_{21}A_{11}^-A_{12}$ and $A_{11} - A_{12}A_{22}^-A_{21}$ appearing in this theorem are commonly referred to as generalized Schur complements.

Theorem 7.8 Let the $m \times m$ matrix A be partitioned as in (7.1), and suppose that A_{11}^- and A_{22}^- are arbitrary generalized inverses of A_{11} and A_{22} . Then

- (a) if $R(A_{21}) \subset R(A_{22})$ or $R(A'_{12}) \subset R(A'_{22})$,

$$|A| = |A_{22}||A_{11} - A_{12}A_{22}^-A_{21}|,$$

- (b) if $R(A_{12}) \subset R(A_{11})$ or $R(A'_{21}) \subset R(A'_{11})$,

$$|A| = |A_{11}||A_{22} - A_{21}A_{11}^-A_{12}|.$$

Proof. It follows from Theorem 5.25 that $A_{22}A_{22}^-A_{21} = A_{21}$ if $R(A_{21}) \subset R(A_{22})$. Consequently,

$$\begin{aligned} A &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \\ &= \begin{bmatrix} I_{m_1} & A_{12} \\ (0) & A_{22} \end{bmatrix} \begin{bmatrix} A_{11} - A_{12}A_{22}^-A_{21} & (0) \\ A_{22}^-A_{21} & I_{m_2} \end{bmatrix}. \end{aligned}$$

Taking the determinant of both sides of this equation and then applying Theorem 7.3 yields the determinantal identity in (a). It also follows from Theorem 5.25 that $A_{12}A_{22}^-A_{22} = A_{12}$ if $R(A'_{12}) \subset R(A'_{22})$. In this case,

$$\begin{aligned} A &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \\ &= \begin{bmatrix} A_{11} - A_{12}A_{22}^-A_{21} & A_{12}A_{22}^- \\ (0) & I_{m_2} \end{bmatrix} \begin{bmatrix} I_{m_1} & (0) \\ A_{21} & A_{22} \end{bmatrix}. \end{aligned}$$

Taking the determinant of both sides of this equation and then applying Theorem 7.3, we again get the determinantal identity in (a). This establishes (a). Part (b) is proven in a similar fashion. \square

It is important to note that if A is a nonnegative definite matrix, then it satisfies the conditions given in (a) and (b) of Theorem 7.8. To see this, recall that if A is nonnegative definite, it can be written as $A = TT'$, where T is $m \times m$. Partitioning T appropriately, we then get

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} = \begin{bmatrix} T_1 \\ T_2 \end{bmatrix} \begin{bmatrix} T_1' & T_2' \end{bmatrix} = \begin{bmatrix} T_1T_1' & T_1T_2' \\ T_2T_1' & T_2T_2' \end{bmatrix}.$$

Since $A_{22} = T_2T_2'$, $R(A_{22}) = R(T_2)$, and because $A_{21} = T_2T_1'$, $R(A_{21}) \subset R(T_2)$, and so $R(A_{21}) \subset R(A_{22})$. Clearly $R(A'_{22}) = R(A_{22})$ and $R(A'_{12}) = R(A_{21})$, so we also have $R(A'_{12}) \subset R(A'_{22})$; that is, the conditions of Theorem 7.8(a) hold. Similarly, it can be shown that the conditions in (b) hold as well.

We will use the normal distribution as an illustration of an application in which a generalized Schur complement arises naturally.

Example 7.4 Suppose that $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where the covariance matrix Ω is positive semidefinite. Since Ω is singular, \mathbf{x} does not have a density function, and so the derivation of the conditional distribution given in Example 7.3 does not apply in this case. Let \mathbf{x} , $\boldsymbol{\mu}$, and Ω be partitioned as in Example 7.3, and again we want to find the conditional distribution of \mathbf{x}_1 given \mathbf{x}_2 . For any $m \times m$ matrix A of constants, we know that $A\mathbf{x} \sim N_m(A\boldsymbol{\mu}, A\Omega A')$, and in particular, if

$$A = \begin{bmatrix} I_{m_1} & -\Omega_{12}\Omega_{22}^- \\ (0) & I_{m_2} \end{bmatrix},$$

we find that

$$A\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 - \Omega_{12}\Omega_{22}^-\mathbf{x}_2 \\ \mathbf{x}_2 \end{bmatrix}$$

has a multivariate normal distribution with mean vector

$$A\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 - \Omega_{12}\Omega_{22}^-\boldsymbol{\mu}_2 \\ \boldsymbol{\mu}_2 \end{bmatrix}.$$

The covariance matrix is

$$\begin{aligned} A\Omega A' &= \begin{bmatrix} I_{m_1} & -\Omega_{12}\Omega_{22}^- \\ (0) & I_{m_2} \end{bmatrix} \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}' & \Omega_{22} \end{bmatrix} \begin{bmatrix} I_{m_1} & (0) \\ -\Omega_{22}^-\Omega_{12}' & I_{m_2} \end{bmatrix} \\ &= \begin{bmatrix} \Omega_{11} - \Omega_{12}\Omega_{22}^-\Omega_{12}' & (0) \\ \Omega_{12}' & \Omega_{22} \end{bmatrix} \begin{bmatrix} I_{m_1} & (0) \\ -\Omega_{22}^-\Omega_{12}' & I_{m_2} \end{bmatrix} \\ &= \begin{bmatrix} \Omega_{11} - \Omega_{12}\Omega_{22}^-\Omega_{12}' & (0) \\ (0) & \Omega_{22} \end{bmatrix}. \end{aligned} \quad (7.10)$$

In simplifying $A\Omega A'$, we have used the fact that $\Omega_{12}\Omega_{22}^-\Omega_{22} = \Omega_{12}$ and the transpose of this identity, $\Omega_{22}\Omega_{22}'\Omega_{12}' = \Omega_{12}'$. As a result of the zero covariances in (7.10), it follows that $\mathbf{x}_1 - \Omega_{12}\Omega_{22}^-\mathbf{x}_2$ is independently distributed of \mathbf{x}_2 , and so its conditional distribution given \mathbf{x}_2 is the same as its unconditional distribution. Consequently, we have

$$\mathbf{x}_1|\mathbf{x}_2 \sim N_{m_1}(\boldsymbol{\mu}_1 + \Omega_{12}\Omega_{22}^-(\mathbf{x}_2 - \boldsymbol{\mu}_2), \Omega_{11} - \Omega_{12}\Omega_{22}^-\Omega_{12}').$$

7.4 RANK

In this section, we wish to find an expression for the rank of A in term of the submatrices given in (7.1). One special case was already given in Theorem 2.9; if

$$A = \begin{bmatrix} A_{11} & (0) \\ (0) & A_{22} \end{bmatrix},$$

then $\text{rank}(A) = \text{rank}(A_{11}) + \text{rank}(A_{22})$. When $A_{12} \neq (0)$ or $A_{21} \neq (0)$, but A_{11} or A_{22} is nonsingular, Theorem 7.9 can be used to determine the rank of A .

Theorem 7.9 Let A be defined as in (7.1). Then

(a) if A_{22} is nonsingular,

$$\text{rank}(A) = \text{rank}(A_{22}) + \text{rank}(A_{11} - A_{12}A_{22}^{-1}A_{21}),$$

(b) if A_{11} is nonsingular,

$$\text{rank}(A) = \text{rank}(A_{11}) + \text{rank}(A_{22} - A_{21}A_{11}^{-1}A_{12}).$$

Proof. To prove part (a), note that because A_{22} is nonsingular, we can write A as

$$\begin{aligned} A &= \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \\ &= \begin{bmatrix} I_{m_1} & A_{12}A_{22}^{-1} \\ (0) & I_{m_2} \end{bmatrix} \begin{bmatrix} A_{11} - A_{12}A_{22}^{-1}A_{21} & (0) \\ (0) & A_{22} \end{bmatrix} \\ &\quad \times \begin{bmatrix} I_{m_1} & (0) \\ A_{22}^{-1}A_{21} & I_{m_2} \end{bmatrix}. \end{aligned}$$

It follows from Theorem 7.4 that the determinant of the matrix

$$\begin{bmatrix} I_{m_1} & A_{12}A_{22}^{-1} \\ (0) & I_{m_2} \end{bmatrix}$$

is 1, so this matrix is nonsingular. Likewise, the matrix

$$\begin{bmatrix} I_{m_1} & (0) \\ A_{22}^{-1}A_{21} & I_{m_2} \end{bmatrix}$$

is nonsingular, so an application of Theorem 1.10 yields

$$\text{rank} \left\{ \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix} \right\} = \text{rank} \left\{ \begin{bmatrix} A_{11} - A_{12}A_{22}^{-1}A_{21} & (0) \\ (0) & A_{22} \end{bmatrix} \right\}.$$

The result now follows from Theorem 2.9. The proof of part (b) is similar. □

Under similar but slightly stronger conditions than those given in Theorem 7.8, we are able to generalize the results of Theorem 7.9 to situations in which both A_{11} and A_{22} are singular matrices.

Theorem 7.10 Let the $m \times m$ matrix A be partitioned as in (7.1), and suppose that A_{11}^- and A_{22}^- are any generalized inverses of A_{11} and A_{22} .

(a) If $R(A_{21}) \subset R(A_{22})$ and $R(A'_{12}) \subset R(A'_{22})$, then

$$\text{rank}(A) = \text{rank}(A_{22}) + \text{rank}(A_{11} - A_{12}A_{22}^-A_{21}).$$

(b) If $R(A_{12}) \subset R(A_{11})$ and $R(A'_{21}) \subset R(A'_{11})$, then

$$\text{rank}(A) = \text{rank}(A_{11}) + \text{rank}(A_{22} - A_{21}A_{11}^-A_{12}).$$

Proof. We will only prove the result in (a) because the proof of (b) is very similar. Since, by Theorem 5.25, $R(A_{21}) \subset R(A_{22})$ and $R(A'_{12}) \subset R(A'_{22})$ imply that $A_{22}A_{22}^-A_{21} = A_{21}$ and $A_{12}A_{22}^-A_{22} = A_{12}$, it is easily verified that

$$C_1AC_2 = C_3, \quad (7.11)$$

where

$$C_1 = \begin{bmatrix} I_{m_1} & -A_{12}A_{22}^- \\ (0) & I_{m_2} \end{bmatrix}, \quad C_2 = \begin{bmatrix} I_{m_1} & (0) \\ -A_{22}^-A_{21} & I_{m_2} \end{bmatrix},$$

and

$$C_3 = \begin{bmatrix} A_{11} - A_{12}A_{22}^-A_{21} & (0) \\ (0) & A_{22} \end{bmatrix}.$$

Applying Theorem 7.4 to C_1 and C_2 , we find that $|C_1| = 1$ and $|C_2| = 1$, and so it follows Theorem 1.10 that

$$\text{rank}(A) = \text{rank}(C_1AC_2) = \text{rank}(C_3).$$

The result now follows from Theorem 2.9. \square

7.5 GENERALIZED INVERSES

In this section, we will present some results for the one condition generalized inverse A^- and the Moore–Penrose inverse A^+ of a matrix A partitioned in the form as given in (7.1). We begin with the generalized inverse A^- . First we will consider a special case.

Theorem 7.11 Let the $m \times m$ matrix A be partitioned as in (7.1). Suppose that $A_{12} = (0)$ and $A_{21} = (0)$ and that A_{11}^- and A_{22}^- are any generalized inverses of A_{11} and A_{22} . Then

$$\begin{bmatrix} A_{11}^- & (0) \\ (0) & A_{22}^- \end{bmatrix}$$

is a generalized inverse of A .

Proof. Denote the matrix given in the theorem by A^- . The result is obtained by verifying through the use of the identities, $A_{11}A_{11}^-A_{11} = A_{11}$ and $A_{22}A_{22}^-A_{22} = A_{22}$, that $AA^-A = A$. \square

Theorem 7.12 gives conditions under which the formulas given in Theorem 7.1 can be used to obtain a generalized inverse of A if we replace the inverses in those formulas by generalized inverses.

Theorem 7.12 Let the $m \times m$ matrix A be partitioned as in (7.1), and suppose that A_{11}^- and A_{22}^- are any generalized inverses of A_{11} and A_{22} .

(a) If $R(A_{21}) \subset R(A_{22})$ and $R(A'_{12}) \subset R(A'_{22})$, then

$$A^- = \begin{bmatrix} B_1^- & -B_1^- A_{12} A_{22}^- \\ -A_{22}^- A_{21} B_1^- & A_{22}^- + A_{22}^- A_{21} B_1^- A_{12} A_{22}^- \end{bmatrix}$$

is a generalized inverse of A , where $B_1 = A_{11} - A_{12} A_{22}^- A_{21}$.

(b) If $R(A_{12}) \subset R(A_{11})$ and $R(A'_{21}) \subset R(A'_{11})$, then

$$A^- = \begin{bmatrix} A_{11}^- + A_{11}^- A_{12} B_2^- A_{21} A_{11}^- & -A_{11}^- A_{12} B_2^- \\ -B_2^- A_{21} A_{11}^- & B_2^- \end{bmatrix}$$

is a generalized inverse of A , where $B_2 = A_{22} - A_{21} A_{11}^- A_{12}$.

Proof. Suppose the conditions in (a) hold. Then equation (7.11) holds and this can be written as

$$C_1 A C_2 = \begin{bmatrix} B_1 & (0) \\ (0) & A_{22} \end{bmatrix},$$

or equivalently,

$$A = C_1^{-1} \begin{bmatrix} B_1 & (0) \\ (0) & A_{22} \end{bmatrix} C_2^{-1},$$

because C_1 and C_2 are nonsingular matrices. As a result of Theorem 5.23(d), a generalized inverse of A can be computed as

$$A^- = C_2 \begin{bmatrix} B_1 & (0) \\ (0) & A_{22} \end{bmatrix}^- C_1.$$

Now using Theorem 7.11, we get

$$\begin{aligned} A^- &= \begin{bmatrix} I_{m_1} & (0) \\ -A_{22}^- A_{21} & I_{m_2} \end{bmatrix} \begin{bmatrix} B_1^- & (0) \\ (0) & A_{22}^- \end{bmatrix} \begin{bmatrix} I_{m_1} & -A_{12} A_{22}^- \\ (0) & I_{m_2} \end{bmatrix} \\ &= \begin{bmatrix} B_1^- & -B_1^- A_{12} A_{22}^- \\ -A_{22}^- A_{21} B_1^- & A_{22}^- + A_{22}^- A_{21} B_1^- A_{12} A_{22}^- \end{bmatrix}, \end{aligned}$$

thereby proving (a). The proof of (b) is similar. \square

Our next result, which is due to Gross (2000), gives an expression for the Moore–Penrose generalized inverse A^+ of the partitioned matrix A when it is nonnegative definite.

Theorem 7.13 Suppose the $m \times m$ nonnegative definite matrix A is partitioned as in (7.1). Then

$$A^+ = \begin{bmatrix} A_{11}^+ + A_{11}^+ A_{12} B^\sim A_{12}' A_{11}^+ & -A_{11}^+ A_{12} B^\sim \\ -B^\sim A_{12}' A_{11}^+ & B^\sim \end{bmatrix} \\ + \begin{bmatrix} -A_{11}^+ (A_{12} Z + Z' A_{12}') A_{11}^+ & A_{11}^+ Z' \\ Z A_{11}^+ & (0) \end{bmatrix},$$

where

$$B^\sim = [Z \ I_{m_2}] \begin{bmatrix} A_{11}^+ + A_{11}^+ A_{12} B^\sim A_{12}' A_{11}^+ & -A_{11}^+ A_{12} B^\sim \\ -B^\sim A_{12}' A_{11}^+ & B^\sim \end{bmatrix} \begin{bmatrix} Z' \\ I_{m_2} \end{bmatrix}, \\ Z = (I_{m_2} - B^+ B) A_{12}' A_{11}^+ \{I_{m_1} + A_{11}^+ A_{12} (I_{m_2} - B^+ B) A_{12}' A_{11}^+\}^{-1}, \\ B = A_{22} - A_{12}' A_{11}^+ A_{12}.$$

Proof. Since A is nonnegative definite, it can be written as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{12}' & A_{22} \end{bmatrix} = \begin{bmatrix} U' \\ V' \end{bmatrix} \begin{bmatrix} U & V \end{bmatrix} = \begin{bmatrix} U'U & U'V \\ V'U & V'V \end{bmatrix},$$

where U and V are $r \times m_1$ and $r \times m_2$ matrices and $r = \text{rank}(A) \leq m$. Denote the Moore–Penrose inverse of A by G so that

$$G = \begin{bmatrix} G_{11} & G_{12} \\ G_{21} & G_{22} \end{bmatrix} = [U \ V]^+ [U \ V]^{+'}, \quad (7.12)$$

where we have used Theorem 5.3(e). Now it follows from Theorem 5.13 that

$$[U \ V]^+ = \begin{bmatrix} U^+ - U^+ V (C^+ + W) \\ C^+ + W \end{bmatrix}, \quad (7.13)$$

where

$$C = (I_r - U U^+) V, \quad W = Z U^+ (I_r - V C^+)$$

and

$$Z = (I_{m_2} - C^+ C) \{I_{m_2} \\ + (I_{m_2} - C^+ C) V' U^{+'} U^+ V (I_{m_2} - C^+ C)\}^{-1} V' U^{+'}.$$

Since

$$U' C = (U' - U' U U^+) V = (U' - U' U^{+'} U^+) V = (U' - U') V = (0),$$

it follows (Problem 5.14) that $C^+U^{+'} = (0)$. Using this and (7.13) in (7.12), we find that

$$\begin{aligned} G_{11} &= (U'U)^+ + U^+VG_{22}V'U^{+'} \\ &\quad - U^+VWU^{+'} - U^+W'V'U^{+'}, \end{aligned} \quad (7.14)$$

$$G_{12} = -U^+VG_{22} + U^+W', \quad (7.15)$$

$$G_{22} = (C^+ + W)(C^+ + W)'. \quad (7.16)$$

Note that

$$\begin{aligned} C'C &= V'(I_r - UU^+)'(I_r - UU^+)V \\ &= V'(I_r - UU^+)V \\ &= V'V - V'U(U'U)^+U'V \\ &= A_{22} - A'_{12}A_{11}^+A_{12} = B. \end{aligned}$$

Thus, using Theorem 5.3(g), we have

$$\begin{aligned} C^+C &= (C'C)^+C'C = B^+B, \\ V'U^{+'} &= V'U(U'U)^+ = A'_{12}A_{11}^+, \\ U^+V &= (U'U)^+U'V = A_{11}^+A_{12}, \end{aligned}$$

so that

$$\begin{aligned} Z &= (I_{m_2} - B^+B)\{I_{m_2} \\ &\quad + (I_{m_2} - B^+B)A'_{12}A_{11}^+A_{11}^+A_{12}(I_{m_2} - B^+B)\}^{-1}A'_{12}A_{11}^+. \end{aligned} \quad (7.17)$$

A simple application of Theorem 1.9 yields

$$\begin{aligned} &\{I_{m_2} + (I_{m_2} - B^+B)A'_{12}A_{11}^+A_{11}^+A_{12}(I_{m_2} - B^+B)\}^{-1} \\ &= I_{m_2} - (I_{m_2} - B^+B)A'_{12}A_{11}^+\{I_{m_1} + A_{11}^+A_{12}(I_{m_2} - B^+B)A'_{12}A_{11}^+\}^{-1} \\ &\quad \times A_{11}^+A_{12}(I_{m_2} - B^+B). \end{aligned} \quad (7.18)$$

Using (7.18) in (7.17) and simplifying, we get

$$Z = (I_{m_2} - B^+B)A'_{12}A_{11}^+\{I_{m_1} + A_{11}^+A_{12}(I_{m_2} - B^+B)A'_{12}A_{11}^+\}^{-1},$$

which is the formula for Z given in the statement of the theorem. By again using $C^+U^{+'} = (0)$, we also find that

$$WU^{+'} = ZU^+U^{+'} = Z(U'U)^+ = ZA_{11}^+, \quad (7.19)$$

$$WC^{+'} = -ZU^+V(C'C)^+ = -ZA_{11}^+A_{12}B^+, \quad (7.20)$$

and

$$\begin{aligned} WW' &= Z(U^+U^{+'} + U^+VC^+C^{+'}V'U^{+'})Z' \\ &= Z(A_{11}^+ + A_{11}^+A_{12}B^+A_{12}'A_{11}^+)Z'. \end{aligned} \quad (7.21)$$

Substituting (7.19), (7.20), and (7.21) in (7.14), (7.15), and (7.16), we then get

$$\begin{aligned} G_{11} &= A_{11}^+ + A_{11}^+A_{12}G_{22}A_{12}'A_{11}^+ - A_{11}^+A_{12}ZA_{11}^+ - A_{11}^+Z'A_{12}'A_{11}^+, \\ G_{12} &= -A_{11}^+A_{12}G_{22} + A_{11}^+Z', \end{aligned}$$

and

$$\begin{aligned} G_{22} &= (C'C)^+ + WC^{+'} + C^+W' + WW' \\ &= B^+ - ZA_{11}^+A_{12}B^+ - B^+A_{12}'A_{11}^+Z' \\ &\quad + Z(A_{11}^+ + A_{11}^+A_{12}B^+A_{12}'A_{11}^+)Z' \\ &= [Z \ I_{m_2}] \begin{bmatrix} A_{11}^+ + A_{11}^+A_{12}B^+A_{12}'A_{11}^+ & -A_{11}^+A_{12}B^+ \\ -B^+A_{12}'A_{11}^+ & B^+ \end{bmatrix} \begin{bmatrix} Z' \\ I_{m_2} \end{bmatrix} \\ &= B^{\sim}, \end{aligned}$$

and so the proof is complete. \square

7.6 EIGENVALUES

In this section, we explore relationships between the eigenvalues of A as given in (7.1) with those of A_{11} and A_{22} . As usual, we will denote the ordered eigenvalues of a matrix such as A using $\lambda_i(A)$; that is, if the eigenvalues of A are real, then they will be identified as $\lambda_1(A) \geq \dots \geq \lambda_m(A)$. Our first result gives bounds on the eigenvalues of A in terms of the eigenvalues of A_{11} and A_{22} when the matrix A is nonnegative definite.

Theorem 7.14 Suppose the $m \times m$ nonnegative definite matrix A is partitioned as in (7.1). Let h and i be integers between 1 and m inclusive. Then

- (a) $\lambda_{h+i-1}(A) \leq \lambda_h(A_{11}) + \lambda_i(A_{22})$, if $h + i \leq m + 1$,
- (b) $\lambda_{h+i-m}(A) \geq \lambda_h(A_{11}) + \lambda_i(A_{22})$, if $h + i \geq m + 1$,

where $\lambda_h(A_{11}) = 0$ if $h > m_1$ and $\lambda_i(A_{22}) = 0$ if $i > m_2$.

Proof. Let $A^{1/2}$ be the symmetric square root matrix of A , and partition it as $A^{1/2} = [F \ G]$, where F is $m \times m_1$ and G is $m \times m_2$. Since

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{12}' & A_{22} \end{bmatrix} = (A^{1/2})'A^{1/2} = \begin{bmatrix} F'F & F'G \\ G'F & G'G \end{bmatrix},$$

we see that $A_{11} = F'F$ and $A_{22} = G'G$. However, we also have

$$A = A^{1/2}(A^{1/2})' = FF' + GG'. \quad (7.22)$$

Since the nonzero eigenvalues of A_{11} and FF' are the same and the nonzero eigenvalues of A_{22} and GG' are the same, the result follows by applying Theorem 3.23 to (7.22). \square

Theorem 7.14 can be used to obtain bounds for sums of eigenvalues of A . For instance, from Theorem 7.14(a), it follows that

$$\lambda_h(A) \leq \lambda_h(A_{11}) + \lambda_1(A_{22})$$

for $h = 1, \dots, m_1$, which immediately leads to

$$\sum_{h=1}^k \lambda_h(A) \leq \sum_{h=1}^k \lambda_h(A_{11}) + k\lambda_1(A_{22}),$$

for $k = 1, \dots, m_1$. The following extension of Theorem 7.14 gives better bounds on the sums of eigenvalues of A than we get from the repeated application of Theorem 7.14 described above.

Theorem 7.15 Suppose the $m \times m$ nonnegative definite matrix A is partitioned as in (7.1), and let i_1, \dots, i_k be distinct integers with $1 \leq i_j \leq m$ for $j = 1, \dots, k$. Then for $k = 1, \dots, m$,

$$\begin{aligned} \sum_{j=1}^k \{\lambda_{i_j}(A_{11}) + \lambda_{m-k+j}(A_{22})\} &\leq \sum_{j=1}^k \lambda_{i_j}(A) \\ &\leq \sum_{j=1}^k \{\lambda_{i_j}(A_{11}) + \lambda_j(A_{22})\}, \end{aligned}$$

where $\lambda_j(A_{11}) = 0$ if $j > m_1$ and $\lambda_j(A_{22}) = 0$ if $j > m_2$.

Proof. The result follows immediately by applying Theorem 3.24 to (7.22). \square

Theorem 7.14 can be used to obtain bounds on differences between eigenvalues of a symmetric matrix A and corresponding eigenvalues of A_{11} or A_{22} . These bounds are useful in obtaining the asymptotic distribution of the eigenvalues of a random symmetric matrix (Eaton and Tyler, 1991).

Theorem 7.16 Let A be an $m \times m$ symmetric matrix partitioned as in (7.1). If $\lambda_{m_1}(A_{11}) > \lambda_1(A_{22})$, then

$$0 \leq \lambda_j(A) - \lambda_j(A_{11}) \leq \lambda_1(A_{12}A'_{12}) / \{\lambda_j(A_{11}) - \lambda_1(A_{22})\}, \quad (7.23)$$

for $j = 1, \dots, m_1$, and

$$\begin{aligned} 0 &\leq \lambda_{m_2-j+1}(A_{22}) - \lambda_{m-j+1}(A) \\ &\leq \lambda_1(A_{12}A'_{12}) / \{\lambda_{m_1}(A_{11}) - \lambda_{m_2-j+1}(A_{22})\}, \end{aligned} \quad (7.24)$$

for $j = 1, \dots, m_2$.

Proof. Since A_{11} is the leading $m_1 \times m_1$ principal submatrix of A , the lower bound in (7.23) follows immediately from Theorem 3.20. Noting that the upper bound in (7.23) holds for A if and only if the upper bound holds for the matrix $A + \alpha I_m$, where α is an arbitrary constant, we may assume without loss of generality that $\lambda_1(A_{22}) = 0$ because we can take $\alpha = -\lambda_1(A_{22})$. In this case, the lower bound in (7.23) implies that $\lambda_j(A) \geq \lambda_j(A_{11}) > 0$, for $j = 1, \dots, m_1$. Let \hat{A} be the matrix obtained by replacing A_{22} in (7.1) by the null matrix so that

$$\hat{A}^2 = \begin{bmatrix} A_{11}^2 + A_{12}A'_{12} & A_{11}A_{12} \\ A'_{12}A_{11} & A'_{12}A_{12} \end{bmatrix}. \quad (7.25)$$

Since $-A_{22}$ is nonnegative definite and $\hat{A} = A - \text{diag}((0), A_{22})$, it follows from Theorem 3.28 that $\lambda_j(\hat{A}) \geq \lambda_j(A)$, for $j = 1, \dots, m$. Now the eigenvalues of \hat{A}^2 are the squares of the eigenvalues of \hat{A} , but we are not assured that the ordered eigenvalues satisfy $\lambda_j(\hat{A}^2) = \lambda_j^2(\hat{A})$ for all j because \hat{A} is not necessarily nonnegative definite and, hence, may have negative eigenvalues. However, we do know that $\lambda_j(\hat{A}) \geq \lambda_j(A) > 0$ for $j = 1, \dots, m_1$, and

$$\lambda_j(\hat{A}^2) \geq \lambda_j^2(\hat{A}) \geq \lambda_j^2(A), \quad (7.26)$$

for $j = 1, \dots, m_1$. Note also that

$$\begin{aligned} \lambda_j(\hat{A}^2) &\leq \lambda_j(A_{11}^2 + A_{12}A'_{12}) + \lambda_1(A'_{12}A_{12}) \\ &\leq \{\lambda_j(A_{11}^2) + \lambda_1(A_{12}A'_{12})\} + \lambda_1(A'_{12}A_{12}) \\ &= \lambda_j^2(A_{11}) + 2\lambda_1(A_{12}A'_{12}), \end{aligned} \quad (7.27)$$

where the first inequality applied Theorem 7.14(a), the second inequality applied Theorem 3.23, whereas the equality follows from the fact that the positive eigenvalues of $A_{12}A'_{12}$ and $A'_{12}A_{12}$ are the same, and $\lambda_j^2(A_{11}) = \lambda_j(A_{11}^2)$ because $\lambda_{m_1}(A_{11}) > 0$. Combining (7.26) and (7.27), we find that for $j = 1, \dots, m_1$,

$$\lambda_j^2(A) \leq \lambda_j^2(A_{11}) + 2\lambda_1(A_{12}A'_{12}),$$

or equivalently,

$$\begin{aligned} \{\lambda_j(A) - \lambda_j(A_{11})\}\{\lambda_j(A) + \lambda_j(A_{11})\} &= \lambda_j^2(A) - \lambda_j^2(A_{11}) \\ &\leq 2\lambda_1(A_{12}A'_{12}). \end{aligned}$$

Thus,

$$\begin{aligned}\lambda_j(A) - \lambda_j(A_{11}) &\leq 2\lambda_1(A_{12}A'_{12})/\{\lambda_j(A) + \lambda_j(A_{11})\} \\ &\leq \lambda_1(A_{12}A'_{12})/\lambda_j(A_{11}),\end{aligned}$$

because $\lambda_j(A) \geq \lambda_j(A_{11})$. This establishes the upper bound in (7.23) because we are assuming that $\lambda_1(A_{22}) = 0$. The inequalities in (7.24) can be obtained by applying those in (7.23) to $-A$. \square

The bounds given in Theorem 7.16 can be improved on; for instance, see Dürnbgen (1995).

In Theorem 7.17, we use Theorem 7.15 to obtain bounds on the difference between a sum of eigenvalues of A and the corresponding sum of eigenvalues of A_{11} .

Theorem 7.17 Let A be an $m \times m$ symmetric matrix partitioned as in (7.1). If $\lambda_{m_1}(A_{11}) > \lambda_1(A_{22})$, then

$$0 \leq \sum_{j=1}^k \{\lambda_j(A) - \lambda_j(A_{11})\} \leq \sum_{j=1}^k \lambda_j(A_{12}A'_{12})/\{\lambda_k(A_{11}) - \lambda_1(A_{22})\}, \quad (7.28)$$

for $k = 1, \dots, m_1$.

Proof. The lower bound follows from Theorem 3.20. As in the proof of Theorem 7.16, we may assume without loss of generality that $\lambda_1(A_{22}) = 0$. Applying Theorem 7.15 to \hat{A}^2 defined in (7.25), we find that

$$\sum_{j=1}^k \lambda_j(\hat{A}^2) \leq \sum_{j=1}^k \lambda_j(A_{11}^2 + A_{12}A'_{12}) + \sum_{j=1}^k \lambda_j(A'_{12}A_{12}). \quad (7.29)$$

An application of Theorem 3.24 yields

$$\sum_{j=1}^k \lambda_j(A_{11}^2 + A_{12}A'_{12}) \leq \sum_{j=1}^k \lambda_j(A_{11}^2) + \sum_{j=1}^k \lambda_j(A_{12}A'_{12}),$$

and when combined with (7.29), we get

$$\sum_{j=1}^k \lambda_j(\hat{A}^2) \leq \sum_{j=1}^k \lambda_j^2(A_{11}) + 2 \sum_{j=1}^k \lambda_j(A_{12}A'_{12}),$$

because the eigenvalues of A_{11} are positive, and the positive eigenvalues of $A_{12}A'_{12}$ and $A'_{12}A_{12}$ are the same. Now using (7.26), we have

$$\sum_{j=1}^k \lambda_j^2(A) \leq \sum_{j=1}^k \lambda_j^2(A_{11}) + 2 \sum_{j=1}^k \lambda_j(A_{12}A'_{12}),$$

or equivalently,

$$\sum_{j=1}^k \{\lambda_j^2(A) - \lambda_j^2(A_{11})\} \leq 2 \sum_{j=1}^k \lambda_j(A_{12}A'_{12}). \quad (7.30)$$

However,

$$\begin{aligned} \sum_{j=1}^k \{\lambda_j^2(A) - \lambda_j^2(A_{11})\} &= \sum_{j=1}^k \{\lambda_j(A) + \lambda_j(A_{11})\} \{\lambda_j(A) - \lambda_j(A_{11})\} \\ &\geq \{\lambda_k(A) + \lambda_k(A_{11})\} \sum_{j=1}^k \{\lambda_j(A) - \lambda_j(A_{11})\} \\ &\geq 2\lambda_k(A_{11}) \sum_{j=1}^k \{\lambda_j(A) - \lambda_j(A_{11})\}. \end{aligned} \quad (7.31)$$

Combining (7.30) and (7.31) leads to

$$\sum_{j=1}^k \{\lambda_j(A) - \lambda_j(A_{11})\} \leq \sum_{j=1}^k \lambda_j(A_{12}A'_{12}) / \lambda_k(A_{11}),$$

which establishes the upper bound in (7.28) because we have assumed that $\lambda_1(A_{22}) = 0$. \square

Our final result compares the eigenvalues of A with those of the Schur complement $A_{11} - A_{12}A_{22}^{-1}A'_{12}$ when A is positive definite.

Theorem 7.18 Suppose A in (7.1) is positive definite, and let $B_1 = A_{11} - A_{12}A_{22}^{-1}A'_{12}$, $B_2 = A_{22} - A'_{12}A_{11}^{-1}A_{12}$, and $C = -B_1^{-1}A_{12}A_{22}^{-1}$. Then if $\lambda_1(B_1) < \lambda_{m_2}(B_2)$,

$$\begin{aligned} 0 &\leq \sum_{j=1}^k \{\lambda_{m_1-j+1}(B_1) - \lambda_{m-j+1}(A)\} \\ &\leq \frac{\lambda_{m_1-k+1}^2(B_1)}{\{\lambda_{m_1-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2)\}} \sum_{j=1}^k \lambda_j(CC'), \end{aligned}$$

for $k = 1, \dots, m_1$.

Proof. Using Theorem 7.1 and the fact that both A_{11} and A_{22} are nonsingular, the inverse of A can be expressed as

$$A^{-1} = \begin{bmatrix} B_1^{-1} & C \\ C' & B_2^{-1} \end{bmatrix}.$$

Applying Theorem 3.20 to A^{-1} , we have for $j = 1, \dots, m_1$, $\lambda_j(B_1^{-1}) \leq \lambda_j(A^{-1})$, so that $\lambda_{m_1-j+1}^{-1}(B_1) \leq \lambda_{m-j+1}^{-1}(A)$ or, equivalently, $\lambda_{m_1-j+1}(B_1) \geq \lambda_{m-j+1}(A)$. This proves the lower bound. Since $\lambda_{m_1}(B_1^{-1}) = \lambda_1^{-1}(B_1) > \lambda_{m_2}^{-1}(B_2) = \lambda_1(B_2^{-1})$, we can apply Theorem 7.17 to A^{-1} , which leads to

$$\begin{aligned} \sum_{j=1}^k \{ \lambda_{m-j+1}^{-1}(A) - \lambda_{m_1-j+1}^{-1}(B_1) \} \\ \leq \sum_{j=1}^k \lambda_j(CC') / \{ \lambda_{m_1-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2) \}. \end{aligned}$$

However,

$$\begin{aligned} \sum_{j=1}^k \{ \lambda_{m-j+1}^{-1}(A) - \lambda_{m_1-j+1}^{-1}(B_1) \} \\ = \sum_{j=1}^k \frac{\lambda_{m_1-j+1}(B_1) - \lambda_{m-j+1}(A)}{\lambda_{m-j+1}(A)\lambda_{m_1-j+1}(B_1)} \\ \geq \lambda_{m-k+1}^{-1}(A)\lambda_{m_1-k+1}^{-1}(B_1) \sum_{j=1}^k \{ \lambda_{m_1-j+1}(B_1) - \lambda_{m-j+1}(A) \}, \end{aligned}$$

so

$$\begin{aligned} \sum_{j=1}^k \{ \lambda_{m_1-j+1}(B_1) - \lambda_{m-j+1}(A) \} \\ \leq \frac{\lambda_{m-k+1}(A)\lambda_{m_1-k+1}(B_1)}{\{ \lambda_{m_1-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2) \}} \sum_{j=1}^k \lambda_j(CC') \\ \leq \frac{\lambda_{m_1-k+1}^2(B_1)}{\{ \lambda_{m_1-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2) \}} \sum_{j=1}^k \lambda_j(CC'), \end{aligned}$$

thereby establishing the upper bound. \square

PROBLEMS

7.1 Consider the $2m \times 2m$ matrix

$$A = \begin{bmatrix} aI_m & bI_m \\ cI_m & dI_m \end{bmatrix},$$

where a, b, c , and d are nonzero scalars.

- (a) Give an expression for the determinant of A .
 (b) For what values of a , b , c , and d will A be nonsingular?
 (c) Find an expression for A^{-1} .

7.2 Let A be of the form

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & (0) \end{bmatrix},$$

where each submatrix is $m \times m$ and the matrices A_{12} and A_{21} are nonsingular. Find an expression for the inverse of A in terms of A_{11} , A_{12} , and A_{21} by applying equations (7.2)–(7.5).

7.3 Generalize Example 7.2 by obtaining the determinant, conditions for non-singularity, and the inverse of the $2m \times 2m$ matrix

$$A = \begin{bmatrix} aI_m & c\mathbf{1}_m\mathbf{1}_m' \\ d\mathbf{1}_m\mathbf{1}_m' & bI_m \end{bmatrix},$$

where a , b , c , and d are nonzero scalars.

7.4 Let the matrix G be given by

$$G = \begin{bmatrix} A & B & C \\ (0) & D & E \\ (0) & (0) & F \end{bmatrix},$$

where each of the matrices A , D , and F is square and nonsingular. Find the inverse of G .

7.5 Use Theorems 7.1 and 7.4 to find the determinant and inverse of the matrix

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

7.6 Suppose

$$V = \begin{bmatrix} A_1 + A_k & A_k & \cdots & A_k \\ A_k & A_2 + A_k & \cdots & A_k \\ \vdots & \vdots & \ddots & \vdots \\ A_k & A_k & \cdots & A_{k-1} + A_k \end{bmatrix},$$

where A_1, \dots, A_k are $m \times m$ positive definite matrices. Show that V^{-1} has the form

$$\begin{bmatrix} A_1^{-1} - A_1^{-1}B^{-1}A_1^{-1} & -A_1^{-1}B^{-1}A_2^{-1} & \cdots & -A_1^{-1}B^{-1}A_{k-1}^{-1} \\ -A_2^{-1}B^{-1}A_1^{-1} & A_2^{-1} - A_2^{-1}B^{-1}A_2^{-1} & \cdots & -A_2^{-1}B^{-1}A_{k-1}^{-1} \\ \vdots & \vdots & \ddots & \vdots \\ -A_{k-1}^{-1}B^{-1}A_1^{-1} & -A_{k-1}^{-1}B^{-1}A_2^{-1} & \cdots & A_{k-1}^{-1} - A_{k-1}^{-1}B^{-1}A_{k-1}^{-1} \end{bmatrix},$$

where $B = \sum_{i=1}^k A_i^{-1}$.

7.7 Let A be an $m \times m$ matrix partitioned as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix},$$

where A_{11} is $m_1 \times m_1$ and $\text{rank}(A) = \text{rank}(A_{11}) = m_1$.

- (a) Show that $A_{22} = A_{21}A_{11}^{-1}A_{12}$.
 (b) Use the result of part (a) to show that

$$B = \begin{bmatrix} A_{11}^{-1} & (0) \\ (0) & (0) \end{bmatrix}$$

is a generalized inverse of A .

- (c) Show that the Moore–Penrose inverse of A is given by

$$A^+ = \begin{bmatrix} A'_{11} \\ A'_{12} \end{bmatrix} C [A'_{11} \ A'_{21}],$$

where $C = (A_{11}A'_{11} + A_{12}A'_{12})^{-1}A_{11}(A'_{11}A_{11} + A'_{21}A_{21})^{-1}$.

7.8 Let A be a symmetric matrix partitioned as in (7.1). Show that A is positive definite if and only if A_{11} and $A_{22} - A_{21}A_{11}^{-1}A_{12}$ are positive definite.

7.9 Let A be an $m \times m$ positive definite matrix, and let B be its inverse. Partition A and B as

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A'_{12} & A_{22} \end{bmatrix}, \quad B = \begin{bmatrix} B_{11} & B_{12} \\ B'_{12} & B_{22} \end{bmatrix},$$

where A_{11} and B_{11} are $m_1 \times m_1$ matrices. Show that the matrix

$$\begin{bmatrix} A_{11} - B_{11}^{-1} & A_{12} \\ A'_{12} & A_{22} \end{bmatrix}$$

is positive semidefinite with rank of $m - m_1$.

7.10 Let A and B be defined as in Theorem 7.1. If A is positive definite, show that $B_{11} - A_{11}^{-1}$ is nonnegative definite.

7.11 Consider the $m \times m$ matrix

$$A = \begin{bmatrix} A_{11} & \mathbf{a} \\ \mathbf{a}' & a_{mm} \end{bmatrix},$$

where the $(m - 1) \times (m - 1)$ matrix A_{11} is positive definite.

- (a) Prove that $|A| \leq a_{mm}|A_{11}|$ with equality if and only if $\mathbf{a} = \mathbf{0}$.
 (b) Generalize the result of part (a) by proving that if a_{11}, \dots, a_{mm} are the diagonal elements of a positive definite matrix A , then $|A| \leq a_{11} \cdots a_{mm}$ with equality if and only if A is a diagonal matrix.

- 7.12** Let A and B be nonsingular matrices, with A being $m \times m$ and B being $n \times n$. If C is $m \times n$, D is $n \times m$, and $A + CBD$ is nonsingular, then an expression for the inverse of $A + CBD$ utilizing the inverse of $B^{-1} + DA^{-1}C$ was given in Theorem 1.9. Show that $A + CBD$ is nonsingular if and only if $B^{-1} + DA^{-1}C$ is nonsingular by applying Theorem 7.4 to the matrix

$$E = \begin{bmatrix} A & C \\ D & -B^{-1} \end{bmatrix}.$$

- 7.13** Let A be defined as in (7.1), and suppose that B is an $m_2 \times m_1$ matrix. Show that

$$\begin{vmatrix} A_{11} & A_{12} \\ A_{21} + BA_{11} & A_{22} + BA_{12} \end{vmatrix} = |A|.$$

- 7.14** Suppose A is partitioned as in (7.1), $m_1 = m_2$, and A is positive definite. Use Theorem 4.17 to show that $|A_{12}|^2 < |A_{11}||A_{22}|$.

- 7.15** Let A be defined as in (7.1), where $m_1 = m_2$ and $A_{11}A_{21} = A_{21}A_{11}$. Show that

$$|A| = |A_{11}A_{22} - A_{21}A_{12}|.$$

- 7.16** Let $\Gamma = (\Gamma_1, \Gamma_2)$ be an $m \times m$ orthogonal matrix with Γ_1 and Γ_2 being $m \times m_1$ and $m \times m_2$, respectively. Show that if A is an $m \times m$ nonsingular matrix, then

$$|\Gamma'_1 A \Gamma_1| = |A| |\Gamma'_2 A^{-1} \Gamma_2|.$$

- 7.17** Let A be an $m \times m$ nonsingular matrix and B be an $m \times m$ matrix having rank 1. By considering a matrix of the form

$$\begin{bmatrix} A & \mathbf{d} \\ \mathbf{c}' & 1 \end{bmatrix},$$

show that

$$|A + B| = \{1 + \text{tr}(A^{-1}B)\}|A|.$$

- 7.18** Let A be an $m \times m$ symmetric matrix, and let A_k be its leading $k \times k$ principal submatrix.

- (a) Show that if $|A_1| > 0, \dots, |A_{m-1}| > 0$ and $|A_m| \geq 0$, then A is nonnegative definite.
- (b) Give an example of a 2×2 symmetric matrix that has both of its leading principal minors being nonnegative, yet the matrix is not nonnegative definite.

- 7.19** Provide the details of the proof of part (b) of Theorem 7.8.

- 7.20** Let A be an $m \times n$ matrix and B be an $n \times m$ matrix. Show that

$$\text{rank}(I_m - AB) = \text{rank}(I_n - BA) + m - n.$$

- 7.21** Show that the conditions given in Theorem 7.8 are not necessary conditions. For instance, find a matrix for which the conditions given in Theorem 7.8(a) do not hold, yet the corresponding determinantal identity does hold.
- 7.22** Let \mathbf{u} and \mathbf{v} be $m \times 1$ vectors and A be an $m \times m$ matrix. Show that if $b \neq 0$, then

$$\text{rank}(A - b^{-1}\mathbf{u}\mathbf{v}') < \text{rank}(A)$$

if and only if vectors \mathbf{x} and \mathbf{y} exist, such that $\mathbf{u} = A\mathbf{y}$, $\mathbf{v} = A'\mathbf{x}$, and $b = \mathbf{x}'A\mathbf{y}$.

- 7.23** Let A be defined as in (7.1), and suppose that the conditions $R(A_{21}) \subset R(A_{22})$ and $R(A'_{12}) \subset R(A'_{22})$ both hold. Show that $A_{11} - A_{12}A_{22}^-A_{21}$ does not depend on the choice of the generalized inverse A_{22}^- . Give an example of a matrix A such that only one of these conditions holds and the Schur complement $A_{11} - A_{12}A_{22}^-A_{21}$ does depend on the choice of A_{22}^- .
- 7.24** Let A be defined as in (7.1), and suppose that the conditions $R(A_{21}) \subset R(A_{22})$ and $R(A'_{12}) \subset R(A'_{22})$ hold. Show that if A is idempotent, then the generalized Schur complement $A_{11} - A_{12}A_{22}^-A_{21}$ is also idempotent.
- 7.25** Suppose A is partitioned as in (7.1) and define $B = A_{22} - A_{21}A_{11}^-A_{12}$, $C = (I_{m_1} - A_{11}A_{11}^-)A_{12}$, $D = A_{21}(I_{m_1} - A_{11}^-A_{11})$, and $E = (I_{m_2} - C^-C) \{(I_{m_2} - DD^-)B(I_{m_2} - C^-C)\}^-(I_{m_2} - DD^-)$. Show that a generalized inverse of A is given by

$$A^- = \begin{bmatrix} A_{11}^- - A_{11}^-A_{12}C^-F - GD^-A_{21}A_{11}^- - GD^-BC^-F & GD^- \\ C^-F & (0) \end{bmatrix} + \begin{bmatrix} GD^-B + A_{11}^-A_{12} \\ -I_{m_2} \end{bmatrix} E [BC^-F + A_{21}A_{11}^- \quad -I_{m_2}],$$

where $F = I_{m_1} - A_{11}A_{11}^-$ and $G = I_{m_1} - A_{11}^-A_{11}$.

- 7.26** Provide the details of the proof of part (b) of Theorem 7.12.
- 7.27** Let A be defined as in (7.1), and let $B_1 = A_{11} - A_{12}A_{22}^-A_{21}$. Show that the matrix

$$\begin{bmatrix} B_1^- & -B_1^-A_{12}A_{22}^- \\ -A_{22}^-A_{21}B_1^- & A_{22}^- + A_{22}^-A_{21}B_1^-A_{12}A_{22}^- \end{bmatrix}$$

is a generalized inverse of A if and only if

- (a) $(I_{m_2} - A_{22}A_{22}^-)A_{21}(I_{m_1} - B_1^-B_1) = (0)$,
 (b) $(I_{m_1} - B_1B_1^-)A_{12}(I_{m_2} - A_{22}^-A_{22}) = (0)$,
 (c) $(I_{m_2} - A_{22}A_{22}^-)A_{21}B_1^-A_{12}(I_{m_2} - A_{22}^-A_{22}) = (0)$.
- 7.28** Let A be defined as in (7.1), and let B_2 be the Schur complement $B_2 = A_{22} - A_{21}A_{11}^+A_{12}$. Consider the matrix

$$C = \begin{bmatrix} A_{11}^+ + A_{11}^+A_{12}B_2^+A_{21}A_{11}^+ & -A_{11}^+A_{12}B_2^+ \\ -B_2^+A_{21}A_{11}^+ & B_2^+ \end{bmatrix}.$$

Show that C is the Moore–Penrose inverse of A if the conditions

- (a) $R(A_{12}) \subset R(A_{11})$ and $R(A'_{21}) \subset R(A'_{11})$,
 (b) $R(A_{21}) \subset R(B_2)$ and $R(A'_{12}) \subset R(B'_2)$
 both hold.

7.29 Consider the nonnegative definite matrix

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

Partition A so that A_{11} is 2×2 . Use Theorem 7.13 to obtain the Moore–Penrose inverse of A .

7.30 Consider the matrices B and B^\sim given in Theorem 7.13. Show that B^\sim is a generalized inverse of B .

7.31 Consider the $m \times m$ matrix

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A'_{12} & (0) \end{bmatrix},$$

where A_{11} is an $m_1 \times m_1$ nonnegative definite matrix and A_{12} is $m_1 \times m_2$. Show that

$$A^+ = \begin{bmatrix} B^+ & A'_{12} - B^+ A_{11} A'_{12} \\ A_{12}^+ - A_{12}^+ A_{11} B^+ & A_{12}^+ A_{11} B^+ A_{11} A_{12}^+ - A_{12}^+ A_{11} A_{12}^+ \end{bmatrix},$$

where $B = C A_{11} C$ and $C = I_{m_1} - A_{12} A_{12}^+ = I_{m_1} - A_{12}^+ A'_{12}$.

7.32 Let A be an $m \times m$ symmetric matrix partitioned as in (7.1). Show that if $\lambda_{m_1}(A_{11}) > \lambda_1(A_{22}) > 0$, then

$$\begin{aligned} 0 &\leq \sum_{j=1}^k \{\lambda_j^2(A) - \lambda_j^2(A_{11})\} \\ &\leq 2 \left\{ 1 + \frac{\lambda_1(A_{22})}{\lambda_k(A_{11}) - \lambda_1(A_{22})} \right\} \sum_{j=1}^k \lambda_j(A_{12} A'_{12}) \end{aligned}$$

for $k = 1, \dots, m_1$.

7.33 Show that, under the conditions of Theorem 7.18,

$$\begin{aligned} 0 &\leq \sum_{j=1}^k \{\lambda_{m_1-j+1}^2(B_1) - \lambda_{m-j+1}^2(A)\} \\ &\leq 2\lambda_{m_1-k+1}^4(B_1) \left\{ 1 + \frac{\lambda_{m_2}^{-1}(B_2)}{\lambda_{m_1-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2)} \right\} \sum_{j=1}^k \lambda_j(CC'), \end{aligned}$$

for $k = 1, \dots, m_1$.

- 7.34** Suppose that A in (7.1) is nonnegative definite with A_{22} being positive definite. Consider the Schur complement $B_1 = A_{11} - A_{12}A_{22}^{-1}A'_{12}$. Show that

$$\lambda_{h+m_2}(A) \leq \lambda_h(B_1) \leq \lambda_h(A),$$

for $h = 1, \dots, m_1$.

- 7.35** Suppose that A in (7.1) is nonnegative definite with A_{22} being positive definite, whereas $\text{rank}(B_1) = r$, where $B_1 = A_{11} - A_{12}A_{22}^{-1}A'_{12}$. Let Q be any $m_1 \times r$ matrix satisfying $Q'Q = I_r$ and $B_1 = Q\Delta Q'$, where Δ is a diagonal matrix with the positive eigenvalues of B_1 as its diagonal elements. Define $B_2 = A_{22} - A'_{12}Q(Q'A_{11}Q)^{-1}Q'A_{12}$ and $\hat{C} = -\Delta^{-1}Q'A_{12}A_{22}^{-1}$. Show that, if $\lambda_1(B_1) < \lambda_{m_2}(B_2)$, then

$$\begin{aligned} 0 &\leq \sum_{j=1}^{m_1-r+k} \{\lambda_{m_1-j+1}(B_1) - \lambda_{m-j+1}(A)\} \\ &\leq \frac{\lambda_{r-k+1}^2(B_1)}{\{\lambda_{r-k+1}^{-1}(B_1) - \lambda_{m_2}^{-1}(B_2)\}} \sum_{j=1}^k \lambda_j(\hat{C}\hat{C}'), \end{aligned}$$

for $k = 1, \dots, r$.

8

SPECIAL MATRICES AND MATRIX OPERATIONS

8.1 INTRODUCTION

In this chapter, we will introduce and develop properties of some special matrix operators. In particular, we will look at two matrix products that differ from the ordinary product of two matrices. One of these products, known as the Hadamard product, simply involves the element-wise multiplication of the matrices. The other matrix product, called the Kronecker product, produces a matrix that in partitioned form has each of its submatrices being equal to an element from the first matrix times the second matrix. Closely related to the Kronecker product is the vec operator, which transforms a matrix to a vector by stacking its columns one underneath another. In many situations, a seemingly complicated matrix expression can be written in a fairly simple form by applying one or more of these matrix operators. In addition to these operators, we will look at some special types of structured matrices that we have not previously discussed and are important in some statistical applications.

8.2 THE KRONECKER PRODUCT

Some matrices possess a special type of structure that permits them to be expressed as a product, commonly referred to as the Kronecker product, of two other matrices.

If A is an $m \times n$ matrix and B is a $p \times q$ matrix, then the Kronecker product of A and B , denoted by $A \otimes B$, is the $mp \times nq$ matrix

$$\begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix}. \quad (8.1)$$

This Kronecker product is more precisely known as the right Kronecker product, and it is the most common definition of the Kronecker product appearing in the literature. However, some authors (for example, Graybill, 1983) define the Kronecker product as the left Kronecker product, which has $B \otimes A$ as the matrix given in (8.1). Throughout this book, any reference to the Kronecker product refers to the right Kronecker product. The special structure of the matrix given in (8.1) leads to simplified formulas for the computation of such things as its inverse, determinant, and eigenvalues. In this section, we will develop some of these formulas as well as some of the more basic properties of the Kronecker product.

Unlike ordinary matrix multiplication, the Kronecker product $A \otimes B$ is defined regardless of the sizes of A and B . However, as with ordinary matrix multiplication, the Kronecker product is not, in general, commutative as is demonstrated in the following example.

Example 8.1 Let A and B be the 1×3 and 2×2 matrices given by

$$A = [0 \ 1 \ 2], \quad B = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}.$$

Then we find that

$$A \otimes B = [0B \ 1B \ 2B] = \begin{bmatrix} 0 & 0 & 1 & 2 & 2 & 4 \\ 0 & 0 & 3 & 4 & 6 & 8 \end{bmatrix},$$

whereas

$$B \otimes A = \begin{bmatrix} 1A & 2A \\ 3A & 4A \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 & 0 & 2 & 4 \\ 0 & 3 & 6 & 0 & 4 & 8 \end{bmatrix}.$$

Some of the basic properties of the Kronecker product, which are easily proven from its definition, are summarized in Theorem 8.1. The proofs are left to the reader as an exercise.

Theorem 8.1 Let A , B , and C be any matrices and \mathbf{a} and \mathbf{b} be any two vectors. Then

- (a) $\alpha \otimes A = A \otimes \alpha = \alpha A$, for any scalar α ,
- (b) $(\alpha A) \otimes (\beta B) = \alpha\beta(A \otimes B)$, for any scalars α and β ,

- (c) $(A \otimes B) \otimes C = A \otimes (B \otimes C)$,
- (d) $(A + B) \otimes C = (A \otimes C) + (B \otimes C)$, if A and B are of the same size,
- (e) $A \otimes (B + C) = (A \otimes B) + (A \otimes C)$, if B and C are of the same size,
- (f) $(A \otimes B)' = A' \otimes B'$,
- (g) $\mathbf{a}\mathbf{b}' = \mathbf{a} \otimes \mathbf{b}' = \mathbf{b}' \otimes \mathbf{a}$.

Although as we have already pointed out, $A \otimes B = B \otimes A$ does not hold in general, we see from Theorem 8.1(g) that this commutative property does hold when A and B' are vectors.

We have a useful property involving the Kronecker product and ordinary matrix multiplication in Theorem 8.2.

Theorem 8.2 Let A , B , C , and D be matrices of sizes $m \times h$, $p \times k$, $h \times n$, and $k \times q$, respectively. Then

$$(A \otimes B)(C \otimes D) = AC \otimes BD. \quad (8.2)$$

Proof. The left-hand side of (8.2) is

$$\begin{bmatrix} a_{11}B & \cdots & a_{1h}B \\ \vdots & & \vdots \\ a_{m1}B & \cdots & a_{mh}B \end{bmatrix} \begin{bmatrix} c_{11}D & \cdots & c_{1n}D \\ \vdots & & \vdots \\ c_{h1}D & \cdots & c_{hn}D \end{bmatrix} = \begin{bmatrix} F_{11} & \cdots & F_{1n} \\ \vdots & & \vdots \\ F_{m1} & \cdots & F_{mn} \end{bmatrix},$$

where

$$F_{ij} = \sum_{l=1}^h a_{il}c_{lj}BD = (A)_{i.}(C)_{.j}BD = (AC)_{ij}BD.$$

The right-hand side of (8.2) is

$$AC \otimes BD = \begin{bmatrix} (AC)_{11}BD & \cdots & (AC)_{1n}BD \\ \vdots & & \vdots \\ (AC)_{m1}BD & \cdots & (AC)_{mn}BD \end{bmatrix},$$

and so the result follows. \square

Our next result demonstrates that the trace of the Kronecker product $A \otimes B$ can be expressed in terms of the trace of A and the trace of B when A and B are square matrices.

Theorem 8.3 Let A be an $m \times m$ matrix and B be a $p \times p$ matrix. Then

$$\text{tr}(A \otimes B) = \text{tr}(A)\text{tr}(B).$$

Proof. Using (8.1) when $n = m$, we see that

$$\operatorname{tr}(A \otimes B) = \sum_{i=1}^m a_{ii} \operatorname{tr}(B) = \left(\sum_{i=1}^m a_{ii} \right) \operatorname{tr}(B) = \operatorname{tr}(A) \operatorname{tr}(B),$$

so that the result holds. \square

Theorem 8.3 gives a simplified expression for the trace of a Kronecker product. There is an analogous result for the determinant of a Kronecker product. However, let us first consider the inverse of $A \otimes B$ and the eigenvalues of $A \otimes B$ when A and B are square matrices.

Theorem 8.4 Let A be an $m \times n$ matrix and B be a $p \times q$ matrix. Then

- (a) $(A \otimes B)^{-1} = A^{-1} \otimes B^{-1}$, if $m = n$, $p = q$, and $A \otimes B$ is nonsingular,
- (b) $(A \otimes B)^+ = A^+ \otimes B^+$,
- (c) $(A \otimes B)^- = A^- \otimes B^-$, for any generalized inverses, A^- and B^- , of A and B .

Proof. Using Theorem 8.2, we find that

$$(A^{-1} \otimes B^{-1})(A \otimes B) = (A^{-1}A \otimes B^{-1}B) = I_m \otimes I_p = I_{mp},$$

so (a) holds. We will leave the verification of (b) and (c) as an exercise for the reader. \square

Theorem 8.5 Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of the $m \times m$ matrix A , and let $\theta_1, \dots, \theta_p$ be the eigenvalues of the $p \times p$ matrix B . Then the set of mp eigenvalues of $A \otimes B$ is given by $\{\lambda_i \theta_j : i = 1, \dots, m; j = 1, \dots, p\}$.

Proof. It follows from Theorem 4.12 that nonsingular matrices P and Q exist, such that

$$P^{-1}AP = T_1, \quad Q^{-1}BQ = T_2,$$

where T_1 and T_2 are upper triangular matrices with the eigenvalues of A and B , respectively, as diagonal elements. The eigenvalues of $A \otimes B$ are the same as those of

$$\begin{aligned} (P \otimes Q)^{-1}(A \otimes B)(P \otimes Q) &= (P^{-1} \otimes Q^{-1})(A \otimes B)(P \otimes Q) \\ &= P^{-1}AP \otimes Q^{-1}BQ = T_1 \otimes T_2, \end{aligned}$$

which must be upper triangular because T_1 and T_2 are upper triangular. The result now follows because the eigenvalues of $T_1 \otimes T_2$ are its diagonal elements, which are clearly given by $\{\lambda_i \theta_j : i = 1, \dots, m; j = 1, \dots, p\}$. \square

A simplified expression for the determinant of $A \otimes B$, when A and B are square matrices, is most easily obtained by using the fact that the determinant of a matrix is given by the product of its eigenvalues.

Theorem 8.6 Let A be an $m \times m$ matrix and B be a $p \times p$ matrix. Then

$$|A \otimes B| = |A|^p |B|^m.$$

Proof. Let $\lambda_1, \dots, \lambda_m$ be the eigenvalues of A , and let $\theta_1, \dots, \theta_p$ be the eigenvalues of B . Then we have

$$|A| = \prod_{i=1}^m \lambda_i, \quad |B| = \prod_{j=1}^p \theta_j,$$

and from Theorem 8.5

$$\begin{aligned} |A \otimes B| &= \prod_{j=1}^p \prod_{i=1}^m \lambda_i \theta_j = \prod_{j=1}^p \theta_j^m \left(\prod_{i=1}^m \lambda_i \right) = \prod_{j=1}^p \theta_j^m |A| \\ &= |A|^p \left(\prod_{j=1}^p \theta_j \right)^m = |A|^p |B|^m, \end{aligned}$$

and so the proof is complete. \square

Our final result on Kronecker products identifies a relationship between $\text{rank}(A \otimes B)$ and $\text{rank}(A)$ and $\text{rank}(B)$.

Theorem 8.7 Let A be an $m \times n$ matrix and B be a $p \times q$ matrix. Then

$$\text{rank}(A \otimes B) = \text{rank}(A) \text{rank}(B).$$

Proof. Our proof uses Theorem 3.12, which states that the rank of a symmetric matrix equals the number of its nonzero eigenvalues. Although $A \otimes B$ as given is not necessarily symmetric, the matrix $(A \otimes B)(A \otimes B)'$, as well as AA' and BB' , is symmetric. Now from Theorem 2.8, we have

$$\text{rank}(A \otimes B) = \text{rank}\{(A \otimes B)(A \otimes B)'\} = \text{rank}(AA' \otimes BB').$$

Since $AA' \otimes BB'$ is symmetric, its rank is given by the number of its nonzero eigenvalues. Now if $\lambda_1, \dots, \lambda_m$ are the eigenvalues of AA' , and $\theta_1, \dots, \theta_p$ are the eigenvalues of BB' then, by Theorem 8.5, the eigenvalues of $AA' \otimes BB'$ are given by $\{\lambda_i \theta_j : i = 1, \dots, m; j = 1, \dots, p\}$. Clearly, the number of nonzero values in this set is the number of nonzero λ_i 's times the number of nonzero θ_j 's. However, because

AA' and BB' are symmetric, the number of nonzero λ_i 's is given by $\text{rank}(AA') = \text{rank}(A)$, and the number of nonzero θ_j 's is given by $\text{rank}(BB') = \text{rank}(B)$. The proof is now complete. \square

Example 8.2 The computations involved in an analysis of variance are sometimes particularly well suited for the use of the Kronecker product. For example, consider the univariate one-way classification model,

$$y_{ij} = \mu + \tau_i + \epsilon_{ij},$$

which was discussed in Example 3.16, Example 6.11, and Example 6.12. Suppose that we have the same number of observations available from each of the k treatments, so that $j = 1, \dots, n$ for each i . In this case, the model may be written as

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where $X = (\mathbf{1}_k \otimes \mathbf{1}_n, I_k \otimes \mathbf{1}_n)$, $\boldsymbol{\beta} = (\mu, \tau_1, \dots, \tau_k)'$, $\mathbf{y} = (\mathbf{y}'_1, \dots, \mathbf{y}'_k)'$, and $\mathbf{y}_i = (y_{i1}, \dots, y_{in})'$. Consequently, a least squares solution for $\boldsymbol{\beta}$ is easily computed as

$$\begin{aligned} \hat{\boldsymbol{\beta}} &= (X'X)^- X'\mathbf{y} = \left\{ \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ I_k \otimes \mathbf{1}'_n \end{bmatrix} \begin{bmatrix} \mathbf{1}_k \otimes \mathbf{1}_n & I_k \otimes \mathbf{1}_n \end{bmatrix} \right\}^- \\ &\quad \times \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ I_k \otimes \mathbf{1}'_n \end{bmatrix} \mathbf{y} \\ &= \begin{bmatrix} nk & n\mathbf{1}'_k \\ n\mathbf{1}_k & nI_k \end{bmatrix}^- \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ I_k \otimes \mathbf{1}'_n \end{bmatrix} \mathbf{y} \\ &= \begin{bmatrix} (nk)^{-1} & \mathbf{0}' \\ \mathbf{0} & n^{-1}(I_k - k^{-1}\mathbf{1}_k\mathbf{1}'_k) \end{bmatrix} \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ I_k \otimes \mathbf{1}'_n \end{bmatrix} \mathbf{y} \\ &= \begin{bmatrix} (nk)^{-1}(\mathbf{1}'_k \otimes \mathbf{1}'_n) \\ n^{-1}(I_k \otimes \mathbf{1}'_n) - (nk)^{-1}(\mathbf{1}_k\mathbf{1}'_k \otimes \mathbf{1}'_n) \end{bmatrix} \mathbf{y}, \end{aligned}$$

This yields $\hat{\mu} = \bar{y}$ and $\hat{\tau}_i = \bar{y}_i - \bar{y}$, where

$$\bar{y} = \frac{1}{nk} \sum_{i=1}^k \sum_{j=1}^n y_{ij}, \quad \bar{y}_i = \frac{1}{n} \sum_{j=1}^n y_{ij}.$$

Note that this solution is not unique because X is not full rank, and hence the solution depends on the choice of the generalized inverse of $X'X$. However, for each i , $\mu + \tau_i$ is estimable and its estimate is given by $\hat{\mu} + \hat{\tau}_i = \bar{y}_i$. In addition, the sum of squared errors for the model is always unique and is given by

$$(\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = \mathbf{y}'(I_{nk} - X(X'X)^- X')\mathbf{y}$$

$$\begin{aligned}
&= \mathbf{y}'(I_{nk} - n^{-1}(I_k \otimes \mathbf{1}_n \mathbf{1}_n'))\mathbf{y} \\
&= \sum_{i=1}^k \mathbf{y}_i'(I_n - n^{-1}\mathbf{1}_n \mathbf{1}_n')\mathbf{y}_i \\
&= \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2.
\end{aligned}$$

Since $\{(\mathbf{1}_k \otimes \mathbf{1}_n)'(\mathbf{1}_k \otimes \mathbf{1}_n)\}^{-1}(\mathbf{1}_k \otimes \mathbf{1}_n)'\mathbf{y} = \bar{y}$, the reduced model

$$y_{ij} = \mu + \epsilon_{ij}$$

has the least squares estimate $\hat{\mu} = \bar{y}$, whereas its sum of squared errors is

$$\{\mathbf{y} - \bar{y}(\mathbf{1}_k \otimes \mathbf{1}_n)\}'\{\mathbf{y} - \bar{y}(\mathbf{1}_k \otimes \mathbf{1}_n)\} = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y})^2.$$

The difference in the sums of squared errors for these two models, the so-called sum of squares for treatments (SST), is then

$$\begin{aligned}
\text{SST} &= \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y})^2 - \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2 \\
&= \sum_{i=1}^k n(\bar{y}_i - \bar{y})^2.
\end{aligned}$$

Example 8.3 In this example, we will illustrate some of the computations involved in the analysis of the two-way classification model with interaction, which is of the form

$$y_{ijk} = \mu + \tau_i + \gamma_j + \eta_{ij} + \epsilon_{ijk},$$

where $i = 1, \dots, a$, $j = 1, \dots, b$, and $k = 1, \dots, n$ (see Problem 6.23). Here μ can be described as an overall effect, whereas τ_i is an effect due to the i th level of factor A, γ_j is an effect due to the j th level of factor B, and η_{ij} is an effect due to the interaction of the i th and j th levels of factors A and B. If we define the parameter vector $\boldsymbol{\beta} = (\mu, \tau_1, \dots, \tau_a, \gamma_1, \dots, \gamma_b, \eta_{11}, \eta_{12}, \dots, \eta_{ab-1}, \eta_{ab})'$ and the response vector $\mathbf{y} = (y_{111}, \dots, y_{11n}, y_{121}, \dots, y_{1bn}, y_{211}, \dots, y_{abn})'$, then the model above can be written as

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where

$$X = (\mathbf{1}_a \otimes \mathbf{1}_b \otimes \mathbf{1}_n, I_a \otimes \mathbf{1}_b \otimes \mathbf{1}_n, \mathbf{1}_a \otimes I_b \otimes \mathbf{1}_n, I_a \otimes I_b \otimes \mathbf{1}_n).$$

Now it is easily verified that the matrix

$$X'X = \begin{bmatrix} abn & bn\mathbf{1}'_a & an\mathbf{1}'_b & n\mathbf{1}'_a \otimes \mathbf{1}'_b \\ bn\mathbf{1}_a & bnI_a & n\mathbf{1}_a \otimes \mathbf{1}'_b & nI_a \otimes \mathbf{1}'_b \\ an\mathbf{1}_b & n\mathbf{1}'_a \otimes \mathbf{1}_b & anI_b & n\mathbf{1}'_a \otimes I_b \\ n\mathbf{1}_a \otimes \mathbf{1}_b & nI_a \otimes \mathbf{1}_b & n\mathbf{1}_a \otimes I_b & nI_a \otimes I_b \end{bmatrix}$$

has as a generalized inverse the matrix

$$\text{diag}((abn)^{-1}, (bn)^{-1}(I_a - a^{-1}\mathbf{1}_a\mathbf{1}'_a), (an)^{-1}(I_b - b^{-1}\mathbf{1}_b\mathbf{1}'_b), C),$$

where

$$\begin{aligned} C = & n^{-1}I_a \otimes I_b - (bn)^{-1}I_a \otimes \mathbf{1}_b\mathbf{1}'_b - (an)^{-1}\mathbf{1}_a\mathbf{1}'_a \otimes I_b \\ & + (abn)^{-1}\mathbf{1}_a\mathbf{1}'_a \otimes \mathbf{1}_b\mathbf{1}'_b. \end{aligned}$$

Using this generalized inverse, we find that a least squares solution for β is given by

$$\hat{\beta} = (X'X)^- X'y = \begin{bmatrix} \bar{y}_{..} \\ \bar{y}_{1.} - \bar{y}_{..} \\ \vdots \\ \bar{y}_{a.} - \bar{y}_{..} \\ \bar{y}_{.1} - \bar{y}_{..} \\ \vdots \\ \bar{y}_{.b} - \bar{y}_{..} \\ \bar{y}_{11} - \bar{y}_{1.} - \bar{y}_{.1} + \bar{y}_{..} \\ \vdots \\ \bar{y}_{ab} - \bar{y}_{a.} - \bar{y}_{.b} + \bar{y}_{..} \end{bmatrix},$$

where

$$\begin{aligned} \bar{y}_{..} &= (abn)^{-1} \sum_{i=1}^a \sum_{j=1}^b \sum_{k=1}^n y_{ijk}, & \bar{y}_{i.} &= (bn)^{-1} \sum_{j=1}^b \sum_{k=1}^n y_{ijk}, \\ \bar{y}_{.j} &= (an)^{-1} \sum_{i=1}^a \sum_{k=1}^n y_{ijk}, & \bar{y}_{ij} &= n^{-1} \sum_{k=1}^n y_{ijk}. \end{aligned}$$

Clearly, $\mu + \tau_i + \gamma_j + \eta_{ij}$ is estimable, and its estimate, which is the fitted value for y_{ijk} , is $\hat{\mu} + \hat{\tau}_i + \hat{\gamma}_j + \hat{\eta}_{ij} = \bar{y}_{ij}$. We will leave the computation of some of the sums of squares associated with the analysis of this model for the reader as an exercise.

8.3 THE DIRECT SUM

The direct sum is a matrix operator that transforms several square matrices into one block diagonal matrix with these matrices appearing as the submatrices along the diagonal. Recall that a block diagonal matrix is of the form

$$\text{diag}(A_1, \dots, A_r) = \begin{bmatrix} A_1 & (0) & \cdots & (0) \\ (0) & A_2 & \cdots & (0) \\ \vdots & \vdots & & \vdots \\ (0) & (0) & \cdots & A_r \end{bmatrix},$$

where A_i is an $m_i \times m_i$ matrix. This block diagonal matrix is said to be the direct sum of the matrices A_1, \dots, A_r and is sometimes written as

$$\text{diag}(A_1, \dots, A_r) = A_1 \oplus \cdots \oplus A_r.$$

Clearly, the commutative property does not hold for the direct sum because, for instance,

$$A_1 \oplus A_2 = \begin{bmatrix} A_1 & (0) \\ (0) & A_2 \end{bmatrix} \neq \begin{bmatrix} A_2 & (0) \\ (0) & A_1 \end{bmatrix} = A_2 \oplus A_1,$$

unless $A_1 = A_2$. Direct sums of a matrix with itself can be expressed as Kronecker products; that is, if $A_1 = \cdots = A_r = A$, then

$$A_1 \oplus \cdots \oplus A_r = A \oplus \cdots \oplus A = I_r \otimes A.$$

Some of the basic properties of the direct sum are summarized in the following theorem. The proofs, which are fairly straightforward, are left to the reader.

Theorem 8.8 Let A_1, \dots, A_r be matrices, where A_i is $m_i \times m_i$. Then

- (a) $\text{tr}(A_1 \oplus \cdots \oplus A_r) = \text{tr}(A_1) + \cdots + \text{tr}(A_r)$,
- (b) $|A_1 \oplus \cdots \oplus A_r| = |A_1| \cdots |A_r|$,
- (c) if each A_i is nonsingular, $A = A_1 \oplus \cdots \oplus A_r$ is also nonsingular and $A^{-1} = A_1^{-1} \oplus \cdots \oplus A_r^{-1}$,
- (d) $\text{rank}(A_1 \oplus \cdots \oplus A_r) = \text{rank}(A_1) + \cdots + \text{rank}(A_r)$,
- (e) if the eigenvalues of A_i are denoted by $\lambda_{i,1}, \dots, \lambda_{i,m_i}$, the eigenvalues of $A_1 \oplus \cdots \oplus A_r$ are given by $\{\lambda_{i,j} : i = 1, \dots, r; j = 1, \dots, m_i\}$.

8.4 THE VEC OPERATOR

There are situations in which it is useful to transform a matrix to a vector that has as its elements the elements of the matrix. One such situation in statistics involves

the study of the distribution of the sample covariance matrix S . It is usually more convenient mathematically in distribution theory to express density functions and moments of jointly distributed random variables in terms of the vector with these random variables as its components. Thus, the distribution of the random matrix S is usually given in terms of the vector formed by stacking columns of S , one underneath the other.

The operator that transforms a matrix to a vector is known as the vec operator. If the $m \times n$ matrix A has \mathbf{a}_i as its i th column, then $\text{vec}(A)$ is the $mn \times 1$ vector given by

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

Example 8.4 If A is the 2×3 matrix given by

$$A = \begin{bmatrix} 2 & 0 & 5 \\ 8 & 1 & 3 \end{bmatrix},$$

then $\text{vec}(A)$ is the 6×1 vector given by

$$\text{vec}(A) = \begin{bmatrix} 2 \\ 8 \\ 0 \\ 1 \\ 5 \\ 3 \end{bmatrix}.$$

In this section, we develop some of the basic algebra associated with this operator. For instance, if \mathbf{a} is $m \times 1$ and \mathbf{b} is $n \times 1$, then \mathbf{ab}' is $m \times n$ and

$$\text{vec}(\mathbf{ab}') = \text{vec}([b_1\mathbf{a}, b_2\mathbf{a}, \dots, b_n\mathbf{a}]) = \begin{bmatrix} b_1\mathbf{a} \\ b_2\mathbf{a} \\ \vdots \\ b_n\mathbf{a} \end{bmatrix} = \mathbf{b} \otimes \mathbf{a}.$$

Theorem 8.9 gives this result and some others that follow directly from the definition of the vec operator.

Theorem 8.9 Let \mathbf{a} and \mathbf{b} be any two vectors, whereas A and B are two matrices of the same size. Then

- (a) $\text{vec}(\mathbf{a}) = \text{vec}(\mathbf{a}') = \mathbf{a}$,
- (b) $\text{vec}(\mathbf{ab}') = \mathbf{b} \otimes \mathbf{a}$,
- (c) $\text{vec}(\alpha A + \beta B) = \alpha \text{vec}(A) + \beta \text{vec}(B)$, where α and β are scalars.

The trace of a product of two matrices can be expressed in terms of the vecs of those two matrices. This result is given in Theorem 8.10.

Theorem 8.10 Let A and B both be $m \times n$ matrices. Then

$$\text{tr}(A'B) = \{\text{vec}(A)\}'\text{vec}(B).$$

Proof. As usual, let $\mathbf{a}_1, \dots, \mathbf{a}_n$ denote the columns of A and $\mathbf{b}_1, \dots, \mathbf{b}_n$ denote the columns of B . Then

$$\begin{aligned} \text{tr}(A'B) &= \sum_{i=1}^n (A'B)_{ii} = \sum_{i=1}^n \mathbf{a}'_i \mathbf{b}_i = [\mathbf{a}'_1, \dots, \mathbf{a}'_n] \begin{bmatrix} \mathbf{b}_1 \\ \vdots \\ \mathbf{b}_n \end{bmatrix} \\ &= \{\text{vec}(A)\}'\text{vec}(B), \end{aligned}$$

as is required. \square

A generalization of Theorem 8.9(b) to the situation involving the vec of the product of three matrices is given in Theorem 8.11.

Theorem 8.11 Let A , B , and C be matrices of sizes $m \times n$, $n \times p$, and $p \times q$, respectively. Then

$$\text{vec}(ABC) = (C' \otimes A) \text{vec}(B).$$

Proof. Note that if $\mathbf{b}_1, \dots, \mathbf{b}_p$ are the columns of B , then B can be written as

$$B = \sum_{i=1}^p \mathbf{b}_i \mathbf{e}'_i,$$

where \mathbf{e}_i is the i th column of I_p . Thus,

$$\begin{aligned} \text{vec}(ABC) &= \text{vec} \left\{ A \left(\sum_{i=1}^p \mathbf{b}_i \mathbf{e}'_i \right) C \right\} = \sum_{i=1}^p \text{vec}(A \mathbf{b}_i \mathbf{e}'_i C) \\ &= \sum_{i=1}^p \text{vec}\{ (A \mathbf{b}_i) (C' \mathbf{e}_i)' \} = \sum_{i=1}^p C' \mathbf{e}_i \otimes A \mathbf{b}_i \\ &= (C' \otimes A) \sum_{i=1}^p (\mathbf{e}_i \otimes \mathbf{b}_i), \end{aligned}$$

where the second to last equality follows from Theorem 8.9(b). By again using Theorem 8.9(b), we find that

$$\sum_{i=1}^p (\mathbf{e}_i \otimes \mathbf{b}_i) = \sum_{i=1}^p \text{vec}(\mathbf{b}_i \mathbf{e}_i') = \text{vec} \left(\sum_{i=1}^p \mathbf{b}_i \mathbf{e}_i' \right) = \text{vec}(B),$$

and so the result follows. \square

Example 8.5 The growth curve model discussed in Example 2.18 models an $m \times n$ response matrix Y as $Y = XBZ + E$, where the $m \times p$ and $q \times n$ matrices X and Z are known, the $p \times q$ matrix B contains unknown parameters, and E is an $m \times n$ matrix of errors. In this example, we will illustrate an alternative way of finding the least squares estimator, \hat{B} , of B by utilizing the vec operator and properties of projections. Letting $\hat{Y} = X\hat{B}Z$ denote the matrix of fitted values, \hat{B} is chosen so that

$$\text{tr}\{(Y - \hat{Y})(Y - \hat{Y})'\} = \{\text{vec}(Y) - \text{vec}(\hat{Y})\}'\{\text{vec}(Y) - \text{vec}(\hat{Y})\}$$

is minimized. But since $\text{vec}(\hat{Y}) = \text{vec}(X\hat{B}Z) = (Z' \otimes X) \text{vec}(\hat{B})$, we see that $\text{vec}(\hat{Y})$ must be in the column space of $(Z' \otimes X)$, and so the required minimum is obtained when $\text{vec}(\hat{Y})$ is the projection of $\text{vec}(Y)$ onto that space. Since

$$\begin{aligned} P_{R(Z' \otimes X)} &= (Z' \otimes X)\{(Z' \otimes X)'(Z' \otimes X)\}^+(Z' \otimes X)' \\ &= Z'(ZZ')^+Z \otimes X(X'X)^+X', \end{aligned}$$

this leads to

$$\begin{aligned} \text{vec}(\hat{Y}) &= P_{R(Z' \otimes X)} \text{vec}(Y) \\ &= \text{vec}(X(X'X)^+X'YZ'(ZZ')^+Z), \end{aligned}$$

or equivalently

$$\hat{Y} = X\hat{B}Z = X(X'X)^+X'YZ'(ZZ')^+Z.$$

When $\text{rank}(X) = p$ and $\text{rank}(Z) = q$, \hat{B} is the unique solution obtained by pre-multiplying this equation by $(X'X)^{-1}X'$ and postmultiplying by $Z'(ZZ')^{-1}$; that is,

$$\hat{B} = (X'X)^{-1}X'YZ'(ZZ')^{-1}.$$

Example 8.6 We return to the multivariate multiple regression model,

$$Y = XB + E,$$

considered in Example 3.15. Here Y and E are $N \times m$, X is $N \times k$, and B is $k \times m$. In Example 3.15, it was established that the least squares estimator of B is $\hat{B} = (X'X)^{-1}X'Y$ by showing that the sum of squared error corresponding to an arbitrary estimator B_0 is equal to the sum of squared errors corresponding to the estimator \hat{B} plus a term that is nonnegative. Using the vec operator and the Kronecker product, we

can establish the same result by using the projection method employed in Example 2.11 to find the least squares estimator of β in the standard multiple regression model. The least squares estimator \hat{B} by definition minimizes

$$\{\text{vec}(Y) - \text{vec}(\hat{Y})\}'\{\text{vec}(Y) - \text{vec}(\hat{Y})\}, \quad (8.3)$$

where $\hat{Y} = X\hat{B}$, and $\text{vec}(\hat{Y}) = \text{vec}(X\hat{B}) = (I_m \otimes X) \text{vec}(\hat{B})$, so $\text{vec}(\hat{Y})$ is a point in the subspace of R^{mN} spanned by the columns of $I_m \otimes X$. Since (8.3) is minimized by selecting as this point the orthogonal projection of $\text{vec}(Y)$ onto this subspace, we have

$$\begin{aligned} \text{vec}(X\hat{B}) &= (I_m \otimes X) \{ (I_m \otimes X)' (I_m \otimes X) \}^{-1} (I_m \otimes X)' \text{vec}(Y) \\ &= \{ I_m \otimes X (X'X)^{-1} X' \} \text{vec}(Y) \\ &= \text{vec} \{ X (X'X)^{-1} X' Y \}, \end{aligned}$$

that is, $X\hat{B} = X(X'X)^{-1}X'Y$. Premultiplying this last equation by $(X'X)^{-1}X'$, we arrive at the desired result.

Theorem 8.10 also can be generalized to a result involving the product of more than two matrices.

Theorem 8.12 Let A , B , C , and D be matrices of sizes $m \times n$, $n \times p$, $p \times q$, and $q \times m$, respectively. Then

$$\text{tr}(ABCD) = \{\text{vec}(A')\}'(D' \otimes B) \text{vec}(C).$$

Proof. Using Theorem 8.10, it follows that

$$\text{tr}(ABCD) = \text{tr}\{A(BCD)\} = \{\text{vec}(A')\}' \text{vec}(BCD).$$

From Theorem 8.11, however, we know that $\text{vec}(BCD) = (D' \otimes B) \text{vec}(C)$, and so the proof is complete. \square

The proofs of the following consequences of Theorem 8.12 are left to the reader as an exercise.

Corollary 8.12.1 Let A and C be matrices of sizes $m \times n$ and $n \times m$, respectively, whereas B and D are $n \times n$. Then

- (a) $\text{tr}(ABC) = \{\text{vec}(A')\}'(I_m \otimes B) \text{vec}(C)$,
- (b) $\text{tr}(AD'BDC) = \{\text{vec}(D)\}'(A'C' \otimes B) \text{vec}(D)$.

Other transformations of a matrix, A , to a vector may be useful when the matrix A has some special structure. One such transformation for an $m \times m$ matrix, denoted by $v(A)$, is defined so as to produce the $m(m+1)/2 \times 1$ vector obtained from $\text{vec}(A)$ by deleting from it all of the elements that are above the diagonal of A . Thus, if A is a lower triangular matrix, $v(A)$ contains all of the elements of A except for the zeros in the upper triangular portion of A . Yet another transformation of the $m \times m$ matrix A to a vector is denoted by $\tilde{v}(A)$ and yields the $m(m-1)/2 \times 1$ vector formed from $v(A)$ by deleting from it all of the diagonal elements of A ; that is, $\tilde{v}(A)$ is the vector obtained by stacking only the portion of the columns of A that are below its diagonal. If A is a skew-symmetric matrix, then A can be reconstructed from $\tilde{v}(A)$ because the diagonal elements of A must be zero, whereas $a_{ji} = -a_{ij}$ if $i \neq j$. The notation $v(A)$ and $\tilde{v}(A)$ we use here corresponds to the notation of Magnus (1988). Others (see, for example, Henderson and Searle, 1979) use the notation $\text{vech}(A)$ and $\text{veck}(A)$. In Section 7, we will discuss some transformations that relate the v and \tilde{v} operators to the vec operator.

Example 8.7 The v and \tilde{v} operators are particularly useful when dealing with covariance and correlation matrices. For instance, suppose that we are interested in the distribution of the sample covariance matrix or the distribution of the sample correlation matrix computed from a sample of observations on three different variables. The resulting sample covariance and correlation matrices would be of the form

$$S = \begin{bmatrix} s_{11} & s_{12} & s_{13} \\ s_{12} & s_{22} & s_{23} \\ s_{13} & s_{23} & s_{33} \end{bmatrix}, \quad R = \begin{bmatrix} 1 & r_{12} & r_{13} \\ r_{12} & 1 & r_{23} \\ r_{13} & r_{23} & 1 \end{bmatrix},$$

so that

$$\begin{aligned} \text{vec}(S) &= (s_{11}, s_{12}, s_{13}, s_{12}, s_{22}, s_{23}, s_{13}, s_{23}, s_{33})', \\ \text{vec}(R) &= (1, r_{12}, r_{13}, r_{12}, 1, r_{23}, r_{13}, r_{23}, 1)'. \end{aligned}$$

Since both S and R are symmetric, there are redundant elements in $\text{vec}(S)$ and $\text{vec}(R)$. The elimination of these elements result in $v(S)$ and $v(R)$ given by

$$\begin{aligned} v(S) &= (s_{11}, s_{12}, s_{13}, s_{22}, s_{23}, s_{33})', \\ v(R) &= (1, r_{12}, r_{13}, 1, r_{23}, 1)'. \end{aligned}$$

Finally, by eliminating the nonrandom 1's from $v(R)$, we obtain

$$\tilde{v}(R) = (r_{12}, r_{13}, r_{23})',$$

which contains all of the random variables in R .

8.5 THE HADAMARD PRODUCT

A matrix operator that is a little more obscure than our other matrix operators, but one that is finding increasing applications in statistics, is known as the Hadamard product. This operator, which we will denote by the symbol \odot , simply performs the element-wise multiplication of two matrices; that is, if A and B are each $m \times n$, then

$$A \odot B = \begin{bmatrix} a_{11}b_{11} & \cdots & a_{1n}b_{1n} \\ \vdots & & \vdots \\ a_{m1}b_{m1} & \cdots & a_{mn}b_{mn} \end{bmatrix}.$$

Clearly, this operation is only defined if the two matrices involved are of the same size.

Example 8.8 If A and B are the 2×3 matrices given by

$$A = \begin{bmatrix} 1 & 4 & 2 \\ 0 & 2 & 3 \end{bmatrix}, \quad B = \begin{bmatrix} 3 & 1 & 3 \\ 6 & 5 & 1 \end{bmatrix},$$

then

$$A \odot B = \begin{bmatrix} 3 & 4 & 6 \\ 0 & 10 & 3 \end{bmatrix}.$$

One of the situations in which the Hadamard product finds application in statistics is in expressions for the covariance structure of certain functions of the sample covariance and sample correlation matrices. We will see examples of this later in Section 11.7. In this section, we will investigate some of the properties of this operator. For a more complete treatment, along with some other examples of applications of the operator in statistics, see Styan (1973) and Horn and Johnson (1991). We begin with some elementary properties that follow directly from the definition of the Hadamard product.

Theorem 8.13 Let A , B , and C be $m \times n$ matrices. Then

- (a) $A \odot B = B \odot A$,
- (b) $(A \odot B) \odot C = A \odot (B \odot C)$,
- (c) $(A + B) \odot C = A \odot C + B \odot C$,
- (d) $(A \odot B)' = A' \odot B'$,
- (e) $A \odot (0) = (0)$,
- (f) $A \odot \mathbf{1}_m \mathbf{1}_n' = A$,
- (g) $A \odot I_m = D_A = \text{diag}(a_{11}, \dots, a_{mm})$ if $m = n$,
- (h) $D(A \odot B) = (DA) \odot B = A \odot (DB)$ and $(A \odot B)E = (AE) \odot B = A \odot (BE)$, if D is an $m \times m$ diagonal matrix and E is an $n \times n$ diagonal matrix,
- (i) $\mathbf{a}\mathbf{b}' \odot \mathbf{c}\mathbf{d}' = (\mathbf{a} \odot \mathbf{c})(\mathbf{b} \odot \mathbf{d})'$, where \mathbf{a} and \mathbf{c} are $m \times 1$ vectors and \mathbf{b} and \mathbf{d} are $n \times 1$ vectors.

We will now show how $A \odot B$ is related to the Kronecker product $A \otimes B$; specifically, $A \odot B$ is a submatrix of $A \otimes B$. To see this, define the $m \times m^2$ matrix Ψ_m as

$$\Psi_m = \sum_{i=1}^m \mathbf{e}_{i,m} (\mathbf{e}_{i,m} \otimes \mathbf{e}_{i,m})',$$

where $\mathbf{e}_{i,m}$ is the i th column of the identity matrix I_m . Note that if A and B are $m \times n$, then $\Psi_m(A \otimes B)\Psi'_n$ forms the $m \times n$ submatrix of the $m^2 \times n^2$ matrix $A \otimes B$, which contains rows $1, m+2, 2m+3, \dots, m^2$ and columns $1, n+2, 2n+3, \dots, n^2$. Taking a closer look at this submatrix, we find that

$$\begin{aligned} \Psi_m(A \otimes B)\Psi'_n &= \sum_{i=1}^m \sum_{j=1}^n \mathbf{e}_{i,m} (\mathbf{e}_{i,m} \otimes \mathbf{e}_{i,m})' (A \otimes B) (\mathbf{e}_{j,n} \otimes \mathbf{e}_{j,n}) \mathbf{e}'_{j,n} \\ &= \sum_{i=1}^m \sum_{j=1}^n \mathbf{e}_{i,m} (\mathbf{e}'_{i,m} A \mathbf{e}_{j,n} \otimes \mathbf{e}'_{i,m} B \mathbf{e}_{j,n}) \mathbf{e}'_{j,n} \\ &= \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ij} \mathbf{e}_{i,m} \mathbf{e}'_{j,n} = A \odot B. \end{aligned}$$

Although the rank of $A \odot B$ is not determined, in general, by the rank of A and the rank of B , we do have the following bound.

Theorem 8.14 Let A and B be $m \times n$ matrices. Then

$$\text{rank}(A \odot B) \leq \text{rank}(A) \text{rank}(B).$$

Proof. Using the identity $\Psi_m(A \otimes B)\Psi'_n = A \odot B$, we get

$$\begin{aligned} \text{rank}(A \odot B) &= \text{rank}(\Psi_m(A \otimes B)\Psi'_n) \leq \text{rank}(A \otimes B) \\ &= \text{rank}(A) \text{rank}(B), \end{aligned}$$

where we have used Theorem 2.8(a) and Theorem 8.7. This completes the proof. \square

Example 8.9 Although Theorem 8.14 gives an upper bound for $\text{rank}(A \odot B)$ in terms of $\text{rank}(A)$ and $\text{rank}(B)$, there is no corresponding lower bound. In other words, it is possible that both A and B have full rank, whereas $A \odot B$ has rank equal to 0. For instance, each of the matrices

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

clearly has rank 3, and yet $A \odot B$ has rank 0 because $A \odot B = (0)$.

Theorem 8.15 shows that a bilinear form in a Hadamard product of two matrices may be written as a trace.

Theorem 8.15 Let A and B be $m \times n$ matrices, and let \mathbf{x} and \mathbf{y} be $m \times 1$ and $n \times 1$ vectors, respectively. Then

- (a) $\mathbf{1}'_m (A \odot B) \mathbf{1}_n = \text{tr}(AB')$,
- (b) $\mathbf{x}'(A \odot B)\mathbf{y} = \text{tr}(D_{\mathbf{x}} A D_{\mathbf{y}} B')$,

where $D_{\mathbf{x}} = \text{diag}(x_1, \dots, x_m)$ and $D_{\mathbf{y}} = \text{diag}(y_1, \dots, y_n)$.

Proof. (a) follows because

$$\begin{aligned} \mathbf{1}'_m (A \odot B) \mathbf{1}_n &= \sum_{i=1}^m \sum_{j=1}^n (A \odot B)_{ij} = \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ij} \\ &= \sum_{i=1}^m (A)_{i.} (B')_{.i} = \sum_{i=1}^m (AB')_{ii} = \text{tr}(AB'). \end{aligned}$$

Note also that $\mathbf{x} = D_{\mathbf{x}} \mathbf{1}_m$ and $\mathbf{y} = D_{\mathbf{y}} \mathbf{1}_n$, so that by using (a) above and Theorem 8.13(h), we find that

$$\begin{aligned} \mathbf{x}'(A \odot B)\mathbf{y} &= \mathbf{1}'_m D_{\mathbf{x}} (A \odot B) D_{\mathbf{y}} \mathbf{1}_n = \mathbf{1}'_m (D_{\mathbf{x}} A \odot B D_{\mathbf{y}}) \mathbf{1}_n \\ &= \text{tr}(D_{\mathbf{x}} A D_{\mathbf{y}} B'), \end{aligned}$$

and this proves (b). □

Example 8.10 A correlation matrix is a nonnegative definite matrix with each diagonal element equal to 1. In Example 3.10, we saw how this constraint on the diagonal elements of a correlation matrix affects the possible choices for the eigenvalues and eigenvectors of the correlation matrix. Let P be an $m \times m$ correlation matrix, and let $P = QDQ'$ be its spectral decomposition; that is, Q is an $m \times m$ orthogonal matrix with orthonormal eigenvectors of P as its columns, and D is a diagonal matrix with the nonnegative eigenvalues of P , d_1, \dots, d_m , as its diagonal elements. Suppose we have a particular orthogonal matrix Q , and we wish to determine the possible choices for the diagonal elements of D so that QDQ' has correlation-matrix structure. The constraint on the diagonal elements of P can be expressed as

$$p_{ii} = (QDQ')_{ii} = (Q)_i. D (Q')_{.i} = \sum_{j=1}^m d_j q_{ij}^2 = 1,$$

for $i = 1, \dots, m$. Note that these m equations can be written as the one matrix equation

$$(Q \odot Q) \mathbf{d} = \mathbf{1}_m,$$

where $\mathbf{d} = (d_1, \dots, d_m)'$. Using the results from Chapter 6, we solve this matrix equation to get

$$\mathbf{d} = \mathbf{1}_m + A\mathbf{b}, \quad (8.4)$$

where \mathbf{b} is an arbitrary $r \times 1$ vector, r is the dimension of the null space of $Q \odot Q$, and A is any $m \times r$ matrix whose columns form a basis for the null space of $Q \odot Q$. Any \mathbf{d} obtained from (8.4) will produce a correlation matrix when used in $P = QDQ'$ as long as each $d_i \geq 0$, and any correlation matrix that has the columns of Q as orthonormal eigenvectors will have its vector of eigenvalues as a solution to (8.4). Note that $\mathbf{1}_m$ is a solution to (8.4) regardless of the choice of Q . This is not surprising because $\mathbf{d} = \mathbf{1}_m$ leads to $P = I_m$ and I_m has the spectral decomposition $I_m = QQ'$, where Q can be any $m \times m$ orthogonal matrix. Also, $\mathbf{d} = \mathbf{1}_m$ is the unique solution to $(Q \odot Q)\mathbf{d} = \mathbf{1}_m$ only if $Q \odot Q$ is nonsingular. In other words, if the correlation matrix $P = QDQ' \neq I_m$, then $Q \odot Q$ must be singular.

The following result can be helpful in determining whether the Hadamard product of two symmetric matrices is nonnegative definite or positive definite.

Theorem 8.16 Let A and B each be an $m \times m$ symmetric matrix. Then

- (a) $A \odot B$ is nonnegative definite if A and B are nonnegative definite,
- (b) $A \odot B$ is positive definite if A and B are positive definite.

Proof. Clearly, if A and B are symmetric, then so is $A \odot B$. Let $B = X\Lambda X'$ be the spectral decomposition of B so that $b_{ij} = \sum \lambda_k x_{ik} x_{jk}$, where $\lambda_k \geq 0$ for all k because B is nonnegative definite. Then we find that for any $m \times 1$ vector \mathbf{y} ,

$$\begin{aligned} \mathbf{y}'(A \odot B)\mathbf{y} &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ij} y_i y_j = \sum_{k=1}^m \left(\sum_{i=1}^m \sum_{j=1}^m \lambda_k (y_i x_{ik}) a_{ij} (y_j x_{jk}) \right) \\ &= \sum_{k=1}^m \lambda_k (\mathbf{y} \odot \mathbf{x}_k)' A (\mathbf{y} \odot \mathbf{x}_k), \end{aligned} \quad (8.5)$$

where \mathbf{x}_k represents the k th column of X . Since A is nonnegative definite, the sum in (8.5) must be nonnegative, and so $A \odot B$ is also nonnegative definite. This proves (a). Now if A is positive definite, then (8.5) will be positive for any $\mathbf{y} \neq \mathbf{0}$ that satisfies $\mathbf{y} \odot \mathbf{x}_k \neq \mathbf{0}$ for at least one k for which $\lambda_k > 0$. However, if B is also positive definite, then $\lambda_k > 0$ for all k , and if \mathbf{y} has its h th component $y_h \neq 0$, then $\mathbf{y} \odot \mathbf{x}_k = \mathbf{0}$ for all k only if the h th row of X has all zeros, which is not possible because X is nonsingular. Consequently, there is no $\mathbf{y} \neq \mathbf{0}$ for which (8.5) equals zero, and so (b) follows. \square

Theorem 8.16(b) gives a sufficient condition for the matrix $A \odot B$ to be positive definite. Example 8.11 demonstrates that this condition is not necessary.

Example 8.11 Consider the 2×2 matrices

$$A = \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}.$$

The matrix B is positive definite because, for instance, $B = VV'$, where

$$V = \begin{bmatrix} 2 & 0 \\ 1 & 1 \end{bmatrix}$$

and $\text{rank}(V) = 2$. Clearly, $A \odot B$ is also positive definite because $A \odot B = B$. However, A is not positive definite because $\text{rank}(A) = 1$.

A sufficient condition for the positive definiteness of $A \odot B$, weaker than that given in Theorem 8.16(b), is given in Theorem 8.17.

Theorem 8.17 Let A and B each be an $m \times m$ symmetric matrix. If B is positive definite and A is nonnegative definite with positive diagonal elements, then $A \odot B$ is positive definite.

Proof. We need to show that for any $\mathbf{x} \neq \mathbf{0}$, $\mathbf{x}'(A \odot B)\mathbf{x} > 0$. Since B is positive definite, a nonsingular matrix T exists, such that $B = TT'$. It follows then from Theorem 8.15(b) that

$$\begin{aligned} \mathbf{x}'(A \odot B)\mathbf{x} &= \text{tr}(D_{\mathbf{x}}AD_{\mathbf{x}}B') = \text{tr}(D_{\mathbf{x}}AD_{\mathbf{x}}TT') \\ &= \text{tr}(T'D_{\mathbf{x}}AD_{\mathbf{x}}T). \end{aligned} \quad (8.6)$$

Since A is nonnegative definite, so is $D_{\mathbf{x}}AD_{\mathbf{x}}$. In addition, if $\mathbf{x} \neq \mathbf{0}$ and A has no diagonal elements equal to zero, then $D_{\mathbf{x}}AD_{\mathbf{x}} \neq (0)$; that is, $D_{\mathbf{x}}AD_{\mathbf{x}}$ has rank of at least one, and so it has at least one positive eigenvalue. Since T is nonsingular, $\text{rank}(D_{\mathbf{x}}AD_{\mathbf{x}}) = \text{rank}(T'D_{\mathbf{x}}AD_{\mathbf{x}}T)$, and so $T'D_{\mathbf{x}}AD_{\mathbf{x}}T$ is also nonnegative definite with at least one positive eigenvalue. The result now follows because (8.6) implies that $\mathbf{x}'(A \odot B)\mathbf{x}$ is the sum of the eigenvalues of $T'D_{\mathbf{x}}AD_{\mathbf{x}}T$. \square

The following result, which gives a relationship between the determinant of a positive definite matrix and its diagonal elements, is commonly known as the Hadamard inequality.

Theorem 8.18 If A is an $m \times m$ positive definite matrix, then

$$|A| \leq \prod_{i=1}^m a_{ii},$$

with equality if and only if A is a diagonal matrix.

Proof. Our proof is by induction. If $m = 2$, then

$$|A| = a_{11}a_{22} - a_{12}^2 \leq a_{11}a_{22},$$

with equality if and only if $a_{12} = 0$, and so the result clearly holds when $m = 2$. For general m , use the cofactor expansion formula for the determinant of A to obtain

$$\begin{aligned} |A| &= a_{11} \begin{vmatrix} a_{22} & a_{23} & \cdots & a_{2m} \\ a_{32} & a_{33} & \cdots & a_{3m} \\ \vdots & \vdots & & \vdots \\ a_{m2} & a_{m3} & \cdots & a_{mm} \end{vmatrix} + \begin{vmatrix} 0 & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{vmatrix} \\ &= a_{11}|A_1| + \begin{vmatrix} 0 & \mathbf{a}' \\ \mathbf{a} & A_1 \end{vmatrix}, \end{aligned} \quad (8.7)$$

where A_1 is the $(m-1) \times (m-1)$ submatrix of A formed by deleting the first row and column of A and $\mathbf{a}' = (a_{12}, \dots, a_{1m})$. Since A is positive definite, A_1 also must be positive definite. Consequently, we can use Theorem 7.4(a) to simplify the second term in the right-hand side of (8.7), leading to the equation

$$|A| = a_{11}|A_1| - \mathbf{a}'A_1^{-1}\mathbf{a}|A_1|.$$

Since A_1 and A_1^{-1} are positive definite, it follows that

$$|A| \leq a_{11}|A_1|,$$

with equality if and only if $\mathbf{a} = \mathbf{0}$. Thus, the result holds for the $m \times m$ matrix A if the result holds for the $(m-1) \times (m-1)$ matrix A_1 , and so our induction proof is complete. \square

Corollary 8.18.1 Let B be an $m \times m$ nonsingular matrix. Then

$$|B|^2 \leq \prod_{i=1}^m \left(\sum_{j=1}^m b_{ij}^2 \right),$$

with equality if and only if the rows of B are orthogonal.

Proof. Since B is nonsingular, the matrix $A = BB'$ is positive definite. Note that

$$|A| = |BB'| = |B||B'| = |B|^2$$

and

$$a_{ii} = (BB')_{ii} = (B)_{i\cdot}(B')_{\cdot i} = (B)_{i\cdot}(B)_{i\cdot}' = \sum_{j=1}^m b_{ij}^2,$$

and so the result follows immediately from Theorem 8.18. \square

Theorem 8.18 also holds for positive semidefinite matrices except that, in this case, A need not be diagonal for equality because one or more of its diagonal elements may equal zero. Likewise, Corollary 8.18.1 holds for singular matrices except for the statement concerning equality.

Hadamard's inequality, as given in Theorem 8.18, can be expressed, with the Hadamard product, as

$$|A| \left(\prod_{i=1}^m 1 \right) \leq |A \odot I_m|, \quad (8.8)$$

where the term $(\prod 1)$ corresponds to the product of the diagonal elements of I_m . Theorem 8.20 will show that the inequality (8.8) holds for other matrices besides the identity. However, first we will need Theorem 8.19.

Theorem 8.19 Let A be an $m \times m$ positive definite matrix, and define

$$A_\alpha = A - \alpha e_1 e_1',$$

where $\alpha = |A|/|A_1|$ and A_1 is the $(m-1) \times (m-1)$ submatrix of A formed by deleting its first row and column. Then A_α is nonnegative definite.

Proof. Let A be partitioned as

$$A = \begin{bmatrix} a_{11} & \mathbf{a}' \\ \mathbf{a} & A_1 \end{bmatrix},$$

and note that because A is positive definite, so is A_1 . Thus, using Theorem 7.4, we find that

$$|A| = \begin{vmatrix} a_{11} & \mathbf{a}' \\ \mathbf{a} & A_1 \end{vmatrix} = |A_1|(a_{11} - \mathbf{a}' A_1^{-1} \mathbf{a}),$$

and so $\alpha = |A|/|A_1| = (a_{11} - \mathbf{a}' A_1^{-1} \mathbf{a})$. Consequently, A_α may be written as

$$\begin{aligned} A_\alpha &= \begin{bmatrix} a_{11} & \mathbf{a}' \\ \mathbf{a} & A_1 \end{bmatrix} - \begin{bmatrix} a_{11} - \mathbf{a}' A_1^{-1} \mathbf{a} & \mathbf{0}' \\ \mathbf{0} & (0) \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}' A_1^{-1} \mathbf{a} & \mathbf{a}' \\ \mathbf{a} & A_1 \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{a}' A_1^{-1} \\ I_{m-1} \end{bmatrix} A_1 [A_1^{-1} \mathbf{a} \quad I_{m-1}]. \end{aligned}$$

Since A_1 is positive definite, an $(m-1) \times (m-1)$ matrix T exists, such that $A_1 = TT'$. If we let $U' = T'[A_1^{-1} \mathbf{a} \quad I_{m-1}]$, then $A_\alpha = UU'$, and so A_α is nonnegative definite. \square

Theorem 8.20 Let A and B be $m \times m$ nonnegative definite matrices. Then

$$|A| \prod_{i=1}^m b_{ii} \leq |A \odot B|.$$

Proof. The result follows immediately if A is singular because $|A| = 0$, whereas $|A \odot B| \geq 0$ is guaranteed by Theorem 8.16. For the case in which A is positive definite, we will prove the result by induction. The result holds when $m = 2$, because in this case

$$\begin{aligned} |A \odot B| &= \begin{vmatrix} a_{11}b_{11} & a_{12}b_{12} \\ a_{12}b_{12} & a_{22}b_{22} \end{vmatrix} = a_{11}a_{22}b_{11}b_{22} - (a_{12}b_{12})^2 \\ &= (a_{11}a_{22} - a_{12}^2)b_{11}b_{22} + a_{12}^2(b_{11}b_{22} - b_{12}^2) \\ &= |A|b_{11}b_{22} + a_{12}^2|B| \geq |A|b_{11}b_{22}. \end{aligned}$$

To prove the result for general m , assume that it holds for $m - 1$, so that

$$|A_1| \prod_{i=2}^m b_{ii} \leq |A_1 \odot B_1|, \quad (8.9)$$

where A_1 and B_1 are the submatrices of A and B formed by deleting their first row and first column. From Theorem 8.19, we know that $(A - \alpha e_1 e_1')$ is nonnegative definite, where $\alpha = |A|/|A_1|$. Thus, by using Theorem 8.16(a), Theorem 8.13(c), and the expansion formula for determinants, we find that

$$\begin{aligned} 0 &\leq |(A - \alpha e_1 e_1') \odot B| = |A \odot B - \alpha e_1 e_1' \odot B| \\ &= |A \odot B - \alpha b_{11} e_1 e_1'| \\ &= |A \odot B| - \alpha b_{11} |(A \odot B)_1|, \end{aligned}$$

where $(A \odot B)_1$ denotes the $(m - 1) \times (m - 1)$ submatrix of $A \odot B$ formed by deleting its first row and column. However, $(A \odot B)_1 = A_1 \odot B_1$ so that the inequality above, along with (8.9) and the identity $\alpha|A_1| = |A|$, implies that

$$|A \odot B| \geq \alpha b_{11} |A_1 \odot B_1| \geq \alpha b_{11} \left(|A_1| \prod_{i=2}^m b_{ii} \right) = |A| \prod_{i=1}^m b_{ii}.$$

The proof is now complete. □

Our final results on Hadamard products involve their eigenvalues. First we obtain bounds for each eigenvalue of the matrix $A \odot B$ when A and B are symmetric.

Theorem 8.21 Let A and B be $m \times m$ symmetric matrices. If A and B are non-negative definite, then the i th largest eigenvalue of $A \odot B$ satisfies

$$\lambda_m(A) \left\{ \min_{1 \leq i \leq m} b_{ii} \right\} \leq \lambda_i(A \odot B) \leq \lambda_1(A) \left\{ \max_{1 \leq i \leq m} b_{ii} \right\}.$$

Proof. Since B is nonnegative definite, an $m \times m$ matrix T exists, such that $B = TT'$. Let \mathbf{t}_j be the j th column of T , whereas t_{ij} denotes the (i, j) th element of T . For any $m \times 1$ vector, $\mathbf{x} \neq \mathbf{0}$, we find that

$$\begin{aligned} \mathbf{x}'(A \odot B)\mathbf{x} &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ij} x_i x_j = \sum_{i=1}^m \sum_{j=1}^m a_{ij} \left(\sum_{h=1}^m t_{ih} t_{jh} \right) x_i x_j \\ &= \sum_{h=1}^m \left(\sum_{i=1}^m \sum_{j=1}^m (x_i t_{ih}) a_{ij} (x_j t_{jh}) \right) = \sum_{h=1}^m (\mathbf{x} \odot \mathbf{t}_h)' A (\mathbf{x} \odot \mathbf{t}_h) \\ &\leq \lambda_1(A) \sum_{h=1}^m (\mathbf{x} \odot \mathbf{t}_h)' (\mathbf{x} \odot \mathbf{t}_h) = \lambda_1(A) \sum_{h=1}^m \sum_{j=1}^m x_j^2 t_{jh}^2 \\ &= \lambda_1(A) \sum_{j=1}^m x_j^2 \left(\sum_{h=1}^m t_{jh}^2 \right) = \lambda_1(A) \sum_{j=1}^m x_j^2 b_{jj} \\ &\leq \lambda_1(A) \left\{ \max_{1 \leq i \leq m} b_{ii} \right\} \mathbf{x}' \mathbf{x}, \end{aligned} \tag{8.10}$$

where the first inequality arises from the relation

$$\lambda_1(A) = \max_{\mathbf{y} \neq \mathbf{0}} \frac{\mathbf{y}' A \mathbf{y}}{\mathbf{y}' \mathbf{y}}$$

given in Theorem 3.16, and the last inequality follows because $\lambda_1(A)$ is nonnegative. Using this same result in Theorem 3.16 applied to $A \odot B$, along with (8.10), we find that for any i , $1 \leq i \leq m$,

$$\lambda_i(A \odot B) \leq \lambda_1(A \odot B) = \max_{\mathbf{x} \neq \mathbf{0}} \frac{\mathbf{x}'(A \odot B)\mathbf{x}}{\mathbf{x}' \mathbf{x}} \leq \lambda_1(A) \left\{ \max_{1 \leq i \leq m} b_{ii} \right\},$$

which is the required upper bound on $\lambda_i(A \odot B)$. By using the identity

$$\lambda_m(A) = \min_{\mathbf{y} \neq \mathbf{0}} \frac{\mathbf{y}' A \mathbf{y}}{\mathbf{y}' \mathbf{y}},$$

the lower bound can be established in a similar fashion. □

The bounds given in Theorem 8.21 can be improved; narrower bounds were obtained by Im (1997). Our final result provides an alternative lower bound for the

smallest eigenvalue of $(A \odot B)$. The derivation of this bound will make use of the following result.

Theorem 8.22 Let A be an $m \times m$ positive definite matrix. Then the matrix $(A \odot A^{-1}) - I_m$ is nonnegative definite.

Proof. Let

$$\sum_{i=1}^m \lambda_i \mathbf{x}_i \mathbf{x}_i'$$

be the spectral decomposition of A so that

$$A^{-1} = \sum_{i=1}^m \lambda_i^{-1} \mathbf{x}_i \mathbf{x}_i'.$$

Then

$$\begin{aligned} (A \odot A^{-1}) - I_m &= (A \odot A^{-1}) - I_m \odot I_m \\ &= \left(\sum_{i=1}^m \lambda_i \mathbf{x}_i \mathbf{x}_i' \odot \sum_{j=1}^m \lambda_j^{-1} \mathbf{x}_j \mathbf{x}_j' \right) \\ &\quad - \left(\sum_{i=1}^m \mathbf{x}_i \mathbf{x}_i' \odot \sum_{j=1}^m \mathbf{x}_j \mathbf{x}_j' \right) \\ &= \sum_{i=1}^m \sum_{j=1}^m (\lambda_i \lambda_j^{-1} - 1) (\mathbf{x}_i \mathbf{x}_i' \odot \mathbf{x}_j \mathbf{x}_j') \\ &= \sum_{i \neq j} (\lambda_i \lambda_j^{-1} - 1) (\mathbf{x}_i \mathbf{x}_i' \odot \mathbf{x}_j \mathbf{x}_j') \\ &= \sum_{i < j} (\lambda_i \lambda_j^{-1} + \lambda_j \lambda_i^{-1} - 2) (\mathbf{x}_i \odot \mathbf{x}_j) (\mathbf{x}_i \odot \mathbf{x}_j)' \\ &= XDX', \end{aligned}$$

where X is the $m \times m(m-1)/2$ matrix having $(\mathbf{x}_i \odot \mathbf{x}_j)$, $i < j$, as its columns, whereas D is the diagonal matrix with its corresponding diagonal elements given by $(\lambda_i \lambda_j^{-1} + \lambda_j \lambda_i^{-1} - 2)$, $i < j$. Since A is positive definite, $\lambda_i > 0$ for all i , and so

$$(\lambda_i \lambda_j^{-1} + \lambda_j \lambda_i^{-1} - 2) = \lambda_i^{-1} \lambda_j^{-1} (\lambda_i - \lambda_j)^2 \geq 0.$$

Thus, D is nonnegative definite, and consequently so is XDX' . □

Theorem 8.23 Let A and B be $m \times m$ nonnegative definite matrices. Then

$$\lambda_m(A \odot B) \geq \lambda_m(AB).$$

Proof. As a result of Theorem 8.16, $A \odot B$ is nonnegative definite, and so the inequality is obvious if either A or B is singular because, in this case, AB will have a zero eigenvalue. Suppose that A and B are positive definite, and let T be any matrix such that $TT' = B$. Note that $T'AT - \lambda_m(AB)I_m$ is nonnegative definite because its i th largest eigenvalue is $\lambda_i(T'AT) - \lambda_m(AB)$, and $\lambda_m(AB) = \lambda_m(T'AT)$. As a result,

$$T^{-1'}(T'AT - \lambda_m(AB)I_m)T^{-1} = A - \lambda_m(AB)B^{-1}$$

is also nonnegative definite. Thus, $(A - \lambda_m(AB)B^{-1}) \odot B$ is nonnegative definite because of Theorem 8.16, whereas $\lambda_m(AB)\{(B^{-1} \odot B) - I_m\}$ is nonnegative definite because of Theorem 8.22, and so the sum of these two matrices, which is given by

$$\begin{aligned} & \{(A - \lambda_m(AB)B^{-1}) \odot B\} + \lambda_m(AB)\{(B^{-1} \odot B) - I_m\} \\ &= A \odot B - \lambda_m(AB)(B^{-1} \odot B) + \lambda_m(AB)(B^{-1} \odot B) - \lambda_m(AB)I_m \\ &= A \odot B - \lambda_m(AB)I_m, \end{aligned}$$

is also nonnegative definite. Consequently, for any \mathbf{x} ,

$$\mathbf{x}'(A \odot B)\mathbf{x} \geq \lambda_m(AB)\mathbf{x}'\mathbf{x},$$

and so the result follows from Theorem 3.16. \square

8.6 THE COMMUTATION MATRIX

An $m \times m$ permutation matrix was defined in Section 1.10 to be any matrix that can be obtained from I_m by permuting its columns. In this section, we discuss a special class of permutation matrices, known as commutation matrices, which are very useful when computing the moments of the multivariate normal and related distributions. We will establish some of the basic properties of commutation matrices. A more complete treatment of this subject can be found in Magnus and Neudecker (1979) and Magnus (1988).

Definition 8.1 Let H_{ij} be the $m \times n$ matrix that has its only nonzero element, a one, in the (i, j) th position. Then the $mn \times mn$ commutation matrix, denoted by K_{mn} , is given by

$$K_{mn} = \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij}). \quad (8.11)$$

The matrix H_{ij} can be conveniently expressed in terms of columns from the identity matrices I_m and I_n . If $\mathbf{e}_{i,m}$ is the i th column of I_m and $\mathbf{e}_{j,n}$ is the j th column of I_n , then $H_{ij} = \mathbf{e}_{i,m}\mathbf{e}'_{j,n}$.

Note that, in general, there is more than one commutation matrix of order mn . For example, for $mn = 6$, we have the four commutation matrices, K_{16} , K_{23} , K_{32} , and K_{61} . Using (8.11), it is easy to verify that $K_{16} = K_{61} = I_6$, whereas

$$K_{23} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix},$$

$$K_{32} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

The fact that $K_{32} = K'_{23}$ is not a coincidence, because this is a general property that follows from the definition of K_{mn} .

Theorem 8.24 The commutation matrix satisfies the properties

- (a) $K_{m1} = K_{1m} = I_m$,
- (b) $K'_{mn} = K_{nm}$,
- (c) $K_{mn}^{-1} = K_{nm}$.

Proof. When H_{ij} is $m \times 1$, then $H_{ij} = e_{i,m}$, and so

$$K_{m1} = \sum_{i=1}^m (e_{i,m} \otimes e'_{i,m}) = I_m = \sum_{i=1}^m (e'_{i,m} \otimes e_{i,m}) = K_{1m},$$

which proves (a). To prove (b), note that

$$K'_{mn} = \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij})' = \sum_{i=1}^m \sum_{j=1}^n (H'_{ij} \otimes H_{ij}) = K_{nm}.$$

Finally, (c) follows because

$$H_{ij}H'_{kl} = e_{i,m}e'_{j,n}e_{l,n}e'_{k,m} = \begin{cases} e_{i,m}e'_{k,m}, & \text{if } j = l, \\ (0), & \text{if } j \neq l, \end{cases}$$

$$H'_{ij}H_{kl} = e_{j,n}e'_{i,m}e_{k,m}e'_{l,n} = \begin{cases} e_{j,n}e'_{l,n}, & \text{if } i = k, \\ (0), & \text{if } i \neq k, \end{cases}$$

and so

$$\begin{aligned}
 K_{mn}K_{nm} &= K_{mn}K'_{mn} = \left\{ \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij}) \right\} \left\{ \sum_{k=1}^m \sum_{l=1}^n (H_{kl} \otimes H'_{kl})' \right\} \\
 &= \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^n (H_{ij} H'_{kl} \otimes H'_{ij} H_{kl}) \\
 &= \sum_{i=1}^m \sum_{j=1}^n (e_{i,m} e'_{i,m} \otimes e_{j,n} e'_{j,n}) \\
 &= \left(\sum_{i=1}^m e_{i,m} e'_{i,m} \right) \otimes \left(\sum_{j=1}^n e_{j,n} e'_{j,n} \right) \\
 &= I_m \otimes I_n = I_{mn}.
 \end{aligned}$$

□

Commutation matrices have important relationships with the vec operator and the Kronecker product. For an $m \times n$ matrix A , $\text{vec}(A)$ and $\text{vec}(A')$ are related because they contain the same elements arranged in a different order; that is, an appropriate reordering of the elements of $\text{vec}(A)$ will produce $\text{vec}(A')$. The commutation matrix K_{mn} is the matrix multiplier that transforms $\text{vec}(A)$ to $\text{vec}(A')$.

Theorem 8.25 For any $m \times n$ matrix A ,

$$K_{mn} \text{vec}(A) = \text{vec}(A').$$

Proof. Clearly, because $a_{ij}H'_{ij}$ is the $n \times m$ matrix whose only nonzero element, a_{ij} , is in the (j, i) th position, we have

$$\begin{aligned}
 A' &= \sum_{i=1}^m \sum_{j=1}^n a_{ij} H'_{ij} = \sum_{i=1}^m \sum_{j=1}^n (e'_{i,m} A e_{j,n}) e_{j,n} e'_{i,m} \\
 &= \sum_{i=1}^m \sum_{j=1}^n e_{j,n} (e'_{i,m} A e_{j,n}) e'_{i,m} = \sum_{i=1}^m \sum_{j=1}^n (e_{j,n} e'_{i,m}) A (e_{j,n} e'_{i,m}) \\
 &= \sum_{i=1}^m \sum_{j=1}^n H'_{ij} A H'_{ij}.
 \end{aligned}$$

Taking the vec of both sides and using Theorem 8.11, we get

$$\text{vec}(A') = \text{vec} \left(\sum_{i=1}^m \sum_{j=1}^n H'_{ij} A H'_{ij} \right) = \sum_{i=1}^m \sum_{j=1}^n \text{vec}(H'_{ij} A H'_{ij})$$

$$= \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij}) \text{vec}(A) = K_{mn} \text{vec}(A),$$

and so the result follows. \square

The term *commutation* arises from the fact that commutation matrices provide the factors that allow a Kronecker product to commute. This property is summarized in Theorem 8.26.

Theorem 8.26 Let A be an $m \times n$ matrix, B be a $p \times q$ matrix, \mathbf{x} be an $m \times 1$ vector, and \mathbf{y} be a $p \times 1$ vector. Then

- (a) $K_{pm}(A \otimes B) = (B \otimes A)K_{qn}$,
- (b) $K_{pm}(A \otimes B)K_{nq} = B \otimes A$,
- (c) $K_{pm}(A \otimes \mathbf{y}) = \mathbf{y} \otimes A$,
- (d) $K_{mp}(\mathbf{y} \otimes A) = A \otimes \mathbf{y}$,
- (e) $K_{pm}(\mathbf{x} \otimes \mathbf{y}) = \mathbf{y} \otimes \mathbf{x}$,
- (f) $\text{tr}\{(B \otimes A)K_{mn}\} = \text{tr}(BA)$, if $p = n$ and $q = m$.

Proof. If X is a $q \times n$ matrix, then by using Theorem 8.11 and Theorem 8.25, we find that

$$\begin{aligned} K_{pm}(A \otimes B) \text{vec}(X) &= K_{pm} \text{vec}(BXA') = \text{vec}\{(BXA')'\} \\ &= \text{vec}(AX'B') = (B \otimes A) \text{vec}(X') \\ &= (B \otimes A)K_{qn} \text{vec}(X). \end{aligned}$$

Thus, if X is chosen so that $\text{vec}(X)$ equals the i th column of I_{qn} , we observe that the i th column of $K_{pm}(A \otimes B)$ must be the same as the i th column of $(B \otimes A)K_{qn}$, so (a) follows. Postmultiplying (a) by K_{nq} and then applying Theorem 8.24(c) yields (b). Properties (c)–(e) follow from (a) and Theorem 8.24(a) because

$$\begin{aligned} K_{pm}(A \otimes \mathbf{y}) &= (\mathbf{y} \otimes A)K_{1n} = \mathbf{y} \otimes A, \\ K_{mp}(\mathbf{y} \otimes A) &= (A \otimes \mathbf{y})K_{n1} = A \otimes \mathbf{y}, \\ K_{pm}(\mathbf{x} \otimes \mathbf{y}) &= (\mathbf{y} \otimes \mathbf{x})K_{11} = \mathbf{y} \otimes \mathbf{x}. \end{aligned}$$

Finally, using the definition of the commutation matrix, we get

$$\begin{aligned} \text{tr}\{(B \otimes A)K_{mn}\} &= \sum_{i=1}^m \sum_{j=1}^n \text{tr}\{(B \otimes A)(H_{ij} \otimes H'_{ij})\} \\ &= \sum_{i=1}^m \sum_{j=1}^n \{\text{tr}(BH_{ij})\} \{\text{tr}(AH'_{ij})\} \end{aligned}$$

$$\begin{aligned}
&= \sum_{i=1}^m \sum_{j=1}^n (e'_{j,n} B e_{i,m}) (e'_{i,m} A e_{j,n}) = \sum_{i=1}^m \sum_{j=1}^n b_{ji} a_{ij} \\
&= \sum_{j=1}^n (B)_{j \cdot} (A)_{\cdot j} = \sum_{j=1}^n (BA)_{jj} = \text{tr}(BA),
\end{aligned}$$

which proves (f). \square

The commutation matrix also can be used to obtain a relationship between the vec of a Kronecker product and the Kronecker product of the corresponding vecs.

Theorem 8.27 Let A be an $m \times n$ matrix and B be a $p \times q$ matrix. Then

$$\text{vec}(A \otimes B) = (I_n \otimes K_{qm} \otimes I_p) \{ \text{vec}(A) \otimes \text{vec}(B) \}.$$

Proof. Our proof follows that given by Magnus (1988). Let $\mathbf{a}_1, \dots, \mathbf{a}_n$ be the columns of A and $\mathbf{b}_1, \dots, \mathbf{b}_q$ be the columns of B . Then, because A and B can be written as

$$A = \sum_{i=1}^n \mathbf{a}_i \mathbf{e}'_{i,n}, \quad B = \sum_{j=1}^q \mathbf{b}_j \mathbf{e}'_{j,q},$$

we have

$$\begin{aligned}
\text{vec}(A \otimes B) &= \sum_{i=1}^n \sum_{j=1}^q \text{vec}(\mathbf{a}_i \mathbf{e}'_{i,n} \otimes \mathbf{b}_j \mathbf{e}'_{j,q}) \\
&= \sum_{i=1}^n \sum_{j=1}^q \text{vec}\{(\mathbf{a}_i \otimes \mathbf{b}_j)(\mathbf{e}'_{i,n} \otimes \mathbf{e}'_{j,q})\} \\
&= \sum_{i=1}^n \sum_{j=1}^q \{(e_{i,n} \otimes e_{j,q}) \otimes (\mathbf{a}_i \otimes \mathbf{b}_j)\} \\
&= \sum_{i=1}^n \sum_{j=1}^q \{e_{i,n} \otimes K_{qm}(\mathbf{a}_i \otimes e_{j,q}) \otimes \mathbf{b}_j\} \\
&= \sum_{i=1}^n \sum_{j=1}^q (I_n \otimes K_{qm} \otimes I_p)(e_{i,n} \otimes \mathbf{a}_i \otimes e_{j,q} \otimes \mathbf{b}_j) \\
&= (I_n \otimes K_{qm} \otimes I_p) \left\{ \sum_{i=1}^n (e_{i,n} \otimes \mathbf{a}_i) \otimes \sum_{j=1}^q (e_{j,q} \otimes \mathbf{b}_j) \right\}
\end{aligned}$$

$$\begin{aligned}
&= (I_n \otimes K_{qm} \otimes I_p) \left\{ \sum_{i=1}^n \text{vec}(\mathbf{a}_i \mathbf{e}'_{i,n}) \otimes \sum_{j=1}^q \text{vec}(\mathbf{b}_j \mathbf{e}'_{j,q}) \right\} \\
&= (I_n \otimes K_{qm} \otimes I_p) \{ \text{vec}(A) \otimes \text{vec}(B) \},
\end{aligned}$$

and so the proof is complete. \square

Our next theorem establishes some results for the special commutation matrix K_{mm} . Corresponding results for the general commutation matrix K_{mn} can be found in Magnus and Neudecker (1979) or Magnus (1988).

Theorem 8.28 The commutation matrix K_{mm} has the eigenvalue $+1$ repeated $\frac{1}{2}m(m+1)$ times and the eigenvalue -1 repeated $\frac{1}{2}m(m-1)$ times. In addition,

$$\text{tr}(K_{mm}) = m \quad \text{and} \quad |K_{mm}| = (-1)^{m(m-1)/2}.$$

Proof. Since K_{mm} is real and symmetric, we know from Theorem 3.9 that its eigenvalues are also real. Further, because K_{mm} is orthogonal, the square of each eigenvalue must be 1, so it has eigenvalues $+1$ and -1 only. Let p be the number of eigenvalues equal to -1 , which implies that $m^2 - p$ is the number of eigenvalues equal to $+1$. Since the trace equals the sum of the eigenvalues, we must have $\text{tr}(K_{mm}) = p(-1) + (m^2 - p)(1) = m^2 - 2p$. However, by using basic properties of the trace, we also find that

$$\begin{aligned}
\text{tr}(K_{mm}) &= \text{tr} \left\{ \sum_{i=1}^m \sum_{j=1}^m (\mathbf{e}_i \mathbf{e}'_j \otimes \mathbf{e}_j \mathbf{e}'_i) \right\} = \sum_{i=1}^m \sum_{j=1}^m \text{tr}(\mathbf{e}_i \mathbf{e}'_j \otimes \mathbf{e}_j \mathbf{e}'_i) \\
&= \sum_{i=1}^m \sum_{j=1}^m \{ \text{tr}(\mathbf{e}_i \mathbf{e}'_j) \} \{ \text{tr}(\mathbf{e}_j \mathbf{e}'_i) \} = \sum_{i=1}^m \sum_{j=1}^m (\mathbf{e}'_i \mathbf{e}_j)^2 \\
&= \sum_{i=1}^m 1 = m.
\end{aligned}$$

Evidently, $m^2 - 2p = m$, so that $p = \frac{1}{2}m(m-1)$ as claimed. Finally, the formula given for the determinant follows directly from the fact that the determinant equals the product of the eigenvalues. \square

Commutation matrices can be used to permute the order in Kronecker products of three or more matrices. For instance, suppose that A is $m \times n$, B is $p \times q$, and C is $r \times s$. Then if $K_{r,mp}$ denotes the commutation matrix K_{rh} , where $h = mp$, it immediately follows from Theorem 8.26(b) that

$$K_{r,mp}(A \otimes B \otimes C)K_{nq,s} = (C \otimes A \otimes B).$$

Theorem 8.29 shows us how higher dimensional commutation matrices, such as $K_{r,mp}$, are related to lower dimensional commutation matrices, such as K_{rm} and K_{rp} .

Theorem 8.29 For any positive integers m, n , and p ,

$$K_{np,m} = (I_n \otimes K_{pm})(K_{nm} \otimes I_p) = (I_p \otimes K_{nm})(K_{pm} \otimes I_n).$$

Proof. Let \mathbf{a} , \mathbf{b} , and \mathbf{c} be arbitrary $m \times 1$, $n \times 1$, and $p \times 1$ vectors, respectively. Then using Theorem 8.26(e), we find that

$$\begin{aligned} K_{np,m}(\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}) &= \mathbf{b} \otimes \mathbf{c} \otimes \mathbf{a} = \mathbf{b} \otimes K_{pm}(\mathbf{a} \otimes \mathbf{c}) \\ &= (I_n \otimes K_{pm})(\mathbf{b} \otimes \mathbf{a} \otimes \mathbf{c}) \\ &= (I_n \otimes K_{pm})\{K_{nm}(\mathbf{a} \otimes \mathbf{b}) \otimes \mathbf{c}\} \\ &= (I_n \otimes K_{pm})(K_{nm} \otimes I_p)(\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}). \end{aligned}$$

This implies that $K_{np,m} = (I_n \otimes K_{pm})(K_{nm} \otimes I_p)$ because \mathbf{a} , \mathbf{b} , and \mathbf{c} are arbitrary. The second part is proven in a similar fashion because clearly $K_{np,m} = K_{pn,m}$ and

$$\begin{aligned} K_{pn,m}(\mathbf{a} \otimes \mathbf{c} \otimes \mathbf{b}) &= \mathbf{c} \otimes \mathbf{b} \otimes \mathbf{a} = \mathbf{c} \otimes K_{nm}(\mathbf{a} \otimes \mathbf{b}) \\ &= (I_p \otimes K_{nm})(\mathbf{c} \otimes \mathbf{a} \otimes \mathbf{b}) \\ &= (I_p \otimes K_{nm})\{K_{pm}(\mathbf{a} \otimes \mathbf{c}) \otimes \mathbf{b}\} \\ &= (I_p \otimes K_{nm})(K_{pm} \otimes I_n)(\mathbf{a} \otimes \mathbf{c} \otimes \mathbf{b}). \end{aligned}$$

□

We will see later that the commutation matrix K_{mm} appears in some important matrix moment formulas through the term $N_m = \frac{1}{2}(I_{m^2} + K_{mm})$. Consequently, we will establish some basic properties of N_m .

Theorem 8.30 Let $N_m = \frac{1}{2}(I_{m^2} + K_{mm})$, and let A and B be $m \times m$ matrices. Then

- (a) $N_m = N'_m = N_m^2$,
- (b) $N_m K_{mm} = N_m = K_{mm} N_m$,
- (c) $N_m \text{vec}(A) = \frac{1}{2} \text{vec}(A + A')$,
- (d) $N_m(A \otimes B)N_m = N_m(B \otimes A)N_m$.

Proof. The symmetry of N_m follows from the symmetry of I_{m^2} and K_{mm} , whereas

$$\begin{aligned} N_m^2 &= \frac{1}{4}(I_{m^2} + K_{mm})^2 = \frac{1}{4}(I_{m^2} + 2K_{mm} + K_{mm}^2) \\ &= \frac{1}{2}(I_{m^2} + K_{mm}) = N_m, \end{aligned}$$

because $K_{mm}^2 = I_{m^2}$ follows from the fact that $K_{mm}^{-1} = K_{mm}$. Similarly, (b) follows from the fact that $K_{mm}^2 = I_{m^2}$. Part (c) is an immediate consequence of

$$I_{m^2} \text{vec}(A) = \text{vec}(A), \quad K_{mm} \text{vec}(A) = \text{vec}(A').$$

Finally, note that by using part (b) and Theorem 8.26(b),

$$N_m(A \otimes B)N_m = N_m K_{mm}(B \otimes A)K_{mm}N_m = N_m(B \otimes A)N_m,$$

which proves (d). □

The proof of our final result will be left to the reader as an exercise.

Theorem 8.31 Let A and B be $m \times m$ matrices, such that $A = BB'$. Then

- (a) $N_m(B \otimes B)N_m = (B \otimes B)N_m = N_m(B \otimes B)$,
- (b) $(B \otimes B)N_m(B' \otimes B') = N_m(A \otimes A)$.

8.7 SOME OTHER MATRICES ASSOCIATED WITH THE VEC OPERATOR

In this section, we introduce several other matrices that, like the commutation matrix, have important relationships with the vec operator. However, each of the matrices we discuss here is useful in working with $\text{vec}(A)$ when the matrix A is square and has some particular structure. A more thorough discussion of this and other related material can be found in Magnus (1988).

When the $m \times m$ matrix A is symmetric, then $\text{vec}(A)$ contains redundant elements because $a_{ij} = a_{ji}$, for $i \neq j$. For this reason, we previously defined $\mathbf{v}(A)$ to be the $m(m+1)/2 \times 1$ vector formed by stacking the columns of the lower triangular portion of A . The matrix that transforms $\mathbf{v}(A)$ into $\text{vec}(A)$ is called the duplication matrix; that is, if we denote this duplication matrix by D_m , then for any $m \times m$ symmetric matrix A ,

$$D_m \mathbf{v}(A) = \text{vec}(A). \tag{8.12}$$

For instance, the duplication matrix D_3 is given by

$$D_3 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

For an explicit expression of the $m^2 \times m(m+1)/2$ duplication matrix D_m , refer to Magnus (1988) or Problem 8.63.

Some properties of the duplication matrix and its Moore–Penrose inverse are summarized in Theorem 8.32.

Theorem 8.32 Let D_m be the $m^2 \times m(m+1)/2$ duplication matrix and D_m^+ be its Moore–Penrose inverse. Then

- (a) $\text{rank}(D_m) = m(m+1)/2$,
- (b) $D_m^+ = (D_m' D_m)^{-1} D_m'$,
- (c) $D_m^+ D_m = I_{m(m+1)/2}$,
- (d) $D_m^+ \text{vec}(A) = \text{v}(A)$ for every $m \times m$ symmetric matrix A .

Proof. Clearly, for every $m(m+1)/2 \times 1$ vector \mathbf{x} , an $m \times m$ symmetric matrix A exists, such that $\mathbf{x} = \text{v}(A)$. However, if for some symmetric A , $D_m \text{v}(A) = \mathbf{0}$, then from the definition of D_m , $\text{vec}(A) = \mathbf{0}$, which then implies that $\text{v}(A) = \mathbf{0}$. Thus, $D_m \mathbf{x} = \mathbf{0}$ only if $\mathbf{x} = \mathbf{0}$, and so D_m has full column rank. Parts (b) and (c) follow immediately from (a) and Theorem 5.3, whereas (d) is obtained by premultiplying (8.12) by D_m^+ and then applying (c). \square

The duplication matrix and its Moore–Penrose inverse have some important relationships with K_{mm} and N_m .

Theorem 8.33 Let D_m be the $m^2 \times m(m+1)/2$ duplication matrix and D_m^+ be its Moore–Penrose inverse. Then

- (a) $K_{mm} D_m = N_m D_m = D_m$,
- (b) $D_m^+ K_{mm} = D_m^+ N_m = D_m^+$,
- (c) $D_m D_m^+ = N_m$.

Proof. For any $m \times m$ symmetric matrix A , it follows that

$$\begin{aligned} K_{mm} D_m \text{v}(A) &= K_{mm} \text{vec}(A) = \text{vec}(A') \\ &= \text{vec}(A) = D_m \text{v}(A). \end{aligned} \quad (8.13)$$

Similarly, we have

$$\begin{aligned} N_m D_m \text{v}(A) &= N_m \text{vec}(A) = \frac{1}{2} \text{vec}(A + A') \\ &= \text{vec}(A) = D_m \text{v}(A). \end{aligned} \quad (8.14)$$

Since $\{\text{v}(A) : A \text{ } m \times m \text{ and } A' = A\}$ is all of $m(m+1)/2$ -dimensional space, (8.13) and (8.14) establish (a). To prove (b), take the transpose of (a), premultiply

all three sides by $(D'_m D_m)^{-1}$, and then apply Theorem 8.32(b). We will prove (c) by showing that for any $m \times m$ matrix A ,

$$D_m D_m^+ \text{vec}(A) = N_m \text{vec}(A).$$

If we define $A_* = \frac{1}{2}(A + A')$, then A_* is symmetric and

$$\begin{aligned} N_m \text{vec}(A) &= \frac{1}{2}(I_{m^2} + K_{mm}) \text{vec}(A) = \frac{1}{2}\{\text{vec}(A) + \text{vec}(A')\} \\ &= \text{vec}(A_*). \end{aligned}$$

Using this result and (b), we find that

$$\begin{aligned} D_m D_m^+ \text{vec}(A) &= D_m D_m^+ N_m \text{vec}(A) = D_m D_m^+ \text{vec}(A_*) \\ &= D_m \mathbf{v}(A_*) = \text{vec}(A_*) = N_m \text{vec}(A), \end{aligned}$$

and so the proof is complete. \square

We know from Theorem 8.32 that $D_m^+ \text{vec}(A) = \mathbf{v}(A)$ if A is an $m \times m$ symmetric matrix. Suppose now that A is not symmetric. What will $D_m^+ \text{vec}(A)$ produce? Let $A_* = \frac{1}{2}(A + A')$, and note that because A_* is symmetric, we must have

$$D_m^+ \text{vec}(A_*) = \mathbf{v}(A_*) = \frac{1}{2}\mathbf{v}(A + A').$$

However,

$$\begin{aligned} D_m^+ \text{vec}(A) - D_m^+ \text{vec}(A_*) &= D_m^+ \text{vec}(A - A_*) = D_m^+ \text{vec}\left\{\frac{1}{2}(A - A')\right\} \\ &= \frac{1}{2}D_m^+ \{\text{vec}(A) - \text{vec}(A')\} \\ &= \frac{1}{2}D_m^+ (I_{m^2} - K_{mm}) \text{vec}(A) \\ &= \frac{1}{2}(D_m^+ - D_m^+) \text{vec}(A) = \mathbf{0}, \end{aligned}$$

where we have used Theorem 8.33(b) in the second to the last step. Thus, $D_m^+ \text{vec}(A)$ is the same as $D_m^+ \text{vec}(A_*)$. This result and the analogous expression for $D_m \mathbf{v}(A)$, the derivation of which we leave as an exercise, are summarized in Theorem 8.34.

Theorem 8.34 Let A be an $m \times m$ matrix. Then

- (a) $D_m^+ \text{vec}(A) = \frac{1}{2}\mathbf{v}(A + A')$,
- (b) $D_m \mathbf{v}(A) = \text{vec}(A_L + A'_L - D_A)$,

where A_L is the lower triangular matrix obtained from A by replacing a_{ij} by 0 if $i < j$, and D_A is the diagonal matrix having the same diagonal elements as A .

We will need Theorem 8.35 in Chapter 9.

Theorem 8.35 If A is an $m \times m$ nonsingular matrix, then $D'_m(A \otimes A)D_m$ is nonsingular and its inverse is given by $D_m^+(A^{-1} \otimes A^{-1})D_m^{+'}$.

Proof. To prove the result, we simply show that the product of the two matrices given in Theorem 8.35 yields $I_{m(m+1)/2}$. Using Theorem 8.31(a), Theorem 8.32(c), and Theorem 8.33(a) and (c), we have

$$\begin{aligned} D'_m(A \otimes A)D_m D_m^+(A^{-1} \otimes A^{-1})D_m^{+'} \\ &= D'_m(A \otimes A)N_m(A^{-1} \otimes A^{-1})D_m^{+'} \\ &= D'_m N_m(A \otimes A)(A^{-1} \otimes A^{-1})D_m^{+'} = D'_m N_m D_m^{+'} \\ &= (N_m D_m)' D_m^{+'} = D'_m D_m^{+'} = (D_m^+ D_m)' = I_{m(m+1)/2}, \end{aligned}$$

and so the result follows. \square

We next consider the situation in which the $m \times m$ matrix A is lower triangular. In this case, the elements of $\text{vec}(A)$ are identical to those of $\text{v}(A)$, except that $\text{vec}(A)$ has some additional zeros. We will denote by L'_m the $m^2 \times m(m+1)/2$ matrix that transforms $\text{v}(A)$ into $\text{vec}(A)$; that is, L'_m satisfies

$$L'_m \text{v}(A) = \text{vec}(A). \quad (8.15)$$

Thus, for instance, for $m = 3$,

$$L'_3 = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}.$$

Note that L'_m can be obtained from D_m by replacing $m(m-1)/2$ of the rows of D_m by rows of zeros. The properties of the matrix L_m in Theorem 8.36 can be proven directly from its definition given in (8.15).

Theorem 8.36 The $m(m+1)/2 \times m^2$ matrix L_m satisfies

- (a) $\text{rank}(L_m) = m(m+1)/2$,
- (b) $L_m L'_m = I_{m(m+1)/2}$,
- (c) $L_m^+ = L'_m$,
- (d) $L_m \text{vec}(A) = v(A)$, for every $m \times m$ matrix A .

Proof. Note that if A is lower triangular, then $\text{vec}(A)' \text{vec}(A) = v(A)' v(A)$, and so (8.15) implies

$$v(A)' L_m L'_m v(A) - v(A)' v(A) = v(A)' (L_m L'_m - I_{m(m+1)/2}) v(A) = 0$$

for all lower triangular matrices A . But this can be true only if (b) holds because $\{v(A) : A \text{ } m \times m \text{ and lower triangular}\} = R^{m(m+1)/2}$. Part (a) follows immediately from (b), as does (c) because $L_m^+ = L'_m (L_m L'_m)^{-1}$. To prove (d), note that every matrix A can be written $A = A_L + A_U$, where A_L is lower triangular and A_U is upper triangular with each diagonal element equal to zero. Clearly,

$$0 = \text{vec}(A_L)' \text{vec}(A_U) = v(A_L)' L_m \text{vec}(A_U),$$

and because, for fixed A_U , this must hold for all choices of the lower triangular matrix A_L , it follows that

$$L_m \text{vec}(A_U) = \mathbf{0}.$$

Thus, using this, along with (8.15), (b), and $v(A_L) = v(A)$, we have

$$\begin{aligned} L_m \text{vec}(A) &= L_m \text{vec}(A_L + A_U) = L_m \text{vec}(A_L) \\ &= L_m L'_m v(A_L) = v(A_L) = v(A), \end{aligned}$$

and so the proof is complete. \square

We see from (d) in Theorem 8.36 that L_m is the matrix that eliminates the zeros in $\text{vec}(A)$ coming from the upper triangular portion of A so as to yield $v(A)$. For this reason, L_m is sometimes referred to as the elimination matrix. Theorem 8.37 gives some relationships between L_m and the matrices D_m and N_m . We will leave the proofs of these results as an exercise for the reader.

Theorem 8.37 The elimination matrix L_m satisfies

- (a) $L_m D_m = I_{m(m+1)/2}$,
- (b) $D_m L_m N_m = N_m$,
- (c) $D_m^+ = L_m N_m$.

The last matrix related to $\text{vec}(A)$ that we will discuss is another sort of elimination matrix. Suppose now that the $m \times m$ matrix A is a strictly lower triangular matrix; that is, it is lower triangular and all of its diagonal elements are zero. In this case, $\tilde{v}(A)$ contains all of the relevant elements of A . We denote by \tilde{L}'_m the $m^2 \times m(m-1)/2$ matrix that transforms $\tilde{v}(A)$ into $\text{vec}(A)$; that is,

$$\tilde{L}'_m \tilde{v}(A) = \text{vec}(A).$$

Thus, for $m = 3$, we have

$$\tilde{L}_3 = \begin{bmatrix} 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \end{bmatrix}.$$

Since \tilde{L}_m is similar to L_m , some of its basic properties parallel those of L_m . For instance, Theorem 8.38 is analogous to Theorem 8.36. The proof of this theorem, which we omit, is similar to that of Theorem 8.36.

Theorem 8.38 The $m(m-1)/2 \times m^2$ matrix \tilde{L}_m satisfies

- (a) $\text{rank}(\tilde{L}_m) = m(m-1)/2$,
- (b) $\tilde{L}_m \tilde{L}'_m = I_{m(m-1)/2}$,
- (c) $\tilde{L}_m^+ = \tilde{L}'_m$,
- (d) $\tilde{L}_m \text{vec}(A) = \tilde{v}(A)$, for every $m \times m$ matrix A .

Theorem 8.39 gives some relationships between \tilde{L}_m , L_m , D_m , K_{mm} , and N_m . The proof is left to the reader as an exercise.

Theorem 8.39 The $m(m-1)/2 \times m^2$ matrix \tilde{L}_m satisfies

- (a) $\tilde{L}_m K_{mm} \tilde{L}'_m = (0)$,
- (b) $\tilde{L}_m K_{mm} L'_m = (0)$,
- (c) $\tilde{L}_m D_m = \tilde{L}_m L'_m$,
- (d) $L'_m L_m \tilde{L}'_m = \tilde{L}'_m$,
- (e) $D_m L_m \tilde{L}'_m = 2N_m \tilde{L}'_m$,
- (f) $\tilde{L}_m L'_m L_m \tilde{L}'_m = I_{m(m-1)/2}$.

8.8 NONNEGATIVE MATRICES

The topic of this section, nonnegative and positive matrices, should not be confused with nonnegative definite and positive definite matrices, which we have discussed earlier on several occasions. An $m \times n$ matrix A is a nonnegative matrix,

indicated by $A \geq (0)$, if each element of A is nonnegative. Similarly, A is a positive matrix, indicated by $A > (0)$, if each element of A is positive. We will write $A \geq B$ and $A > B$ to mean that $A - B \geq (0)$ and $A - B > (0)$, respectively. Any matrix A can be transformed to a nonnegative matrix by replacing each of its elements by its absolute value, which will be denoted by $\text{abs}(A)$; that is, if A is an $m \times n$ matrix, then $\text{abs}(A)$ is also an $m \times n$ matrix with (i, j) th element given by $|a_{ij}|$. We will investigate some of the properties of nonnegative square matrices as well as indicate some of their applications in stochastic processes. For a more exhaustive coverage of this topic, the reader is referred to the texts on nonnegative matrices by Berman and Plemmons (1994), Minc (1988), and Seneta (2006), as well as the books by Gantmacher (1959) and Horn and Johnson (2013). Most of the proofs that we present here follow along the lines of the derivations, based on matrix norms, given in Horn and Johnson (2013).

We begin with some results regarding the spectral radius of nonnegative and positive matrices.

Theorem 8.40 Let A be an $m \times m$ matrix and \mathbf{x} be an $m \times 1$ vector. If $A \geq (0)$ and $\mathbf{x} > \mathbf{0}$, then

$$\min_{1 \leq i \leq m} \sum_{j=1}^m a_{ij} \leq \rho(A) \leq \max_{1 \leq i \leq m} \sum_{j=1}^m a_{ij}, \quad (8.16)$$

$$\min_{1 \leq i \leq m} x_i^{-1} \sum_{j=1}^m a_{ij} x_j \leq \rho(A) \leq \max_{1 \leq i \leq m} x_i^{-1} \sum_{j=1}^m a_{ij} x_j, \quad (8.17)$$

with similar inequalities holding when minimizing and maximizing over columns instead of rows.

Proof. Let

$$\alpha = \min_{1 \leq i \leq m} \sum_{j=1}^m a_{ij},$$

and define the $m \times m$ matrix B to have (i, h) th element

$$b_{ih} = \alpha a_{ih} \left(\sum_{j=1}^m a_{ij} \right)^{-1}$$

if $\alpha > 0$ and $b_{ih} = 0$ if $\alpha = 0$. Note that $\|B\|_\infty = \alpha$ and $b_{ih} \leq a_{ih}$, so that $A \geq B$. Clearly, it follows that for any positive integer k , $A^k \geq B^k$, which implies that $\|A^k\|_\infty \geq \|B^k\|_\infty$ or, equivalently,

$$\{\|A^k\|_\infty\}^{1/k} \geq \{\|B^k\|_\infty\}^{1/k}.$$

Taking the limit as $k \rightarrow \infty$, it follows from Theorem 4.28 that $\rho(A) \geq \rho(B)$. However, this proves the lower bound in (8.16) because $\rho(B) = \alpha$ follows from the fact that $\rho(B) \geq \alpha$ because

$$B\mathbf{1}_m = \alpha\mathbf{1}_m$$

and $\rho(B) \leq \|B\|_\infty = \alpha$ due to Theorem 4.23. The upper bound is proven in a similar fashion using

$$\alpha = \max_{1 \leq i \leq m} \sum_{j=1}^m a_{ij}.$$

The bounds in (8.17) follow directly from those in (8.16) because if we define the matrix $C = D_x^{-1}AD_x$, then $C \geq (0)$, $\rho(C) = \rho(A)$, and $c_{ij} = a_{ij}x_i^{-1}x_j$. \square

Theorem 8.41 Let A be an $m \times m$ positive matrix. Then $\rho(A)$ is positive and is an eigenvalue of A . In addition, a positive eigenvector of A corresponding to the eigenvalue $\rho(A)$ exists.

Proof. $\rho(A) > 0$ follows immediately from Theorem 8.40 because A is positive. By the definition of $\rho(A)$, there exists an eigenvalue of A , λ , such that $|\lambda| = \rho(A)$. Let \mathbf{x} be an eigenvector of A corresponding to λ so that $A\mathbf{x} = \lambda\mathbf{x}$. Note that

$$\begin{aligned} \rho(A) \operatorname{abs}(\mathbf{x}) &= |\lambda| \operatorname{abs}(\mathbf{x}) = \operatorname{abs}(\lambda\mathbf{x}) = \operatorname{abs}(A\mathbf{x}) \\ &\leq \operatorname{abs}(A) \operatorname{abs}(\mathbf{x}) = A \operatorname{abs}(\mathbf{x}), \end{aligned}$$

where the inequality clearly follows from the fact that

$$\left| \sum_{j=1}^m a_{ij}x_j \right| \leq \sum_{j=1}^m |a_{ij}| |x_j|$$

for each i . Thus, the vector $\mathbf{y} = A \operatorname{abs}(\mathbf{x}) - \rho(A) \operatorname{abs}(\mathbf{x})$ is nonnegative. The vector $\mathbf{z} = A \operatorname{abs}(\mathbf{x})$ is positive because A is positive and the eigenvector \mathbf{x} must be a nonnull vector. Now if we assume that $\mathbf{y} \neq \mathbf{0}$, then, again, because A is positive, we have

$$\mathbf{0} < A\mathbf{y} = A\mathbf{z} - \rho(A)\mathbf{z},$$

or simply $A\mathbf{z} > \rho(A)\mathbf{z}$. Premultiplying this inequality by D_z^{-1} , we get

$$D_z^{-1}A\mathbf{z} > \rho(A)\mathbf{1}_m$$

or, in other words,

$$z_i^{-1} \sum_{j=1}^m a_{ij}z_j > \rho(A)$$

holds for each i . However, using Theorem 8.40 this implies that $\rho(A) > \rho(A)$. Thus, we must have $\mathbf{y} = \mathbf{0}$. This yields $A \operatorname{abs}(\mathbf{x}) = \rho(A) \operatorname{abs}(\mathbf{x})$, so that $\operatorname{abs}(\mathbf{x})$ is an eigenvector corresponding to $\rho(A)$, and from this we get $\operatorname{abs}(\mathbf{x}) = \rho(A)^{-1} A \operatorname{abs}(\mathbf{x})$, which shows that $\operatorname{abs}(\mathbf{x})$ is positive because $\rho(A) > 0$ and $A \operatorname{abs}(\mathbf{x}) > \mathbf{0}$. This completes the proof. \square

An immediate consequence of the proof of Theorem 8.41 is Corollary 8.41.1.

Corollary 8.41.1 Let A be an $m \times m$ positive matrix, and suppose that λ is an eigenvalue of A satisfying $|\lambda| = \rho(A)$. If \mathbf{x} is any eigenvector corresponding to λ , then

$$A \operatorname{abs}(\mathbf{x}) = \rho(A) \operatorname{abs}(\mathbf{x}).$$

Before determining the dimensionality of the eigenspace associated with the eigenvalue $\rho(A)$, we need Theorem 8.42.

Theorem 8.42 Let \mathbf{x} be an eigenvector corresponding to the eigenvalue λ of the $m \times m$ positive matrix A . If $|\lambda| = \rho(A)$, then some angle θ exists, such that $e^{-i\theta} \mathbf{x} > \mathbf{0}$.

Proof. Note that

$$\operatorname{abs}(A\mathbf{x}) = \operatorname{abs}(\lambda\mathbf{x}) = \rho(A) \operatorname{abs}(\mathbf{x}), \quad (8.18)$$

whereas it follows from Corollary 8.41.1 that

$$A \operatorname{abs}(\mathbf{x}) = \rho(A) \operatorname{abs}(\mathbf{x}). \quad (8.19)$$

Now by using (8.18) and (8.19), we find that

$$\begin{aligned} \rho(A)|x_j| &= |\lambda||x_j| = |\lambda x_j| = \left| \sum_{k=1}^m a_{jk} x_k \right| \leq \sum_{k=1}^m |a_{jk}| |x_k| \\ &= \sum_{k=1}^m a_{jk} |x_k| = \rho(A)|x_j| \end{aligned}$$

holds for each j . Evidently

$$\left| \sum_{k=1}^m a_{jk} x_k \right| = \sum_{k=1}^m a_{jk} |x_k|,$$

and this implies that the possibly complex numbers $a_{jk} x_k = r_k e^{i\theta_k} = r_k (\cos \theta_k + i \sin \theta_k)$, for $k = 1, \dots, m$, have identical angles; that is, some angle θ exists, such

that each $a_{jk}x_k$, for $k = 1, \dots, m$ can be written in the form $a_{jk}x_k = r_k e^{i\theta} = r_k (\cos \theta + i \sin \theta)$. In this case, $e^{-i\theta} a_{jk}x_k = r_k > 0$, which implies that $e^{-i\theta} x_k > 0$ because $a_{jk} > 0$. \square

Theorem 8.43 not only indicates that the eigenspace corresponding to $\rho(A)$ has dimension one, but also that $\rho(A)$ is the only eigenvalue of A having modulus equal to $\rho(A)$.

Theorem 8.43 If A is an $m \times m$ positive matrix, then the dimension of the eigenspace corresponding to the eigenvalue $\rho(A)$ is one. Further, if λ is an eigenvalue of A and $\lambda \neq \rho(A)$, then $|\lambda| < \rho(A)$.

Proof. The first statement will be proven by showing that if \mathbf{u} and \mathbf{v} are nonnull vectors satisfying $A\mathbf{u} = \rho(A)\mathbf{u}$ and $A\mathbf{v} = \rho(A)\mathbf{v}$, then some scalar c exists, such that $\mathbf{v} = c\mathbf{u}$. Now from Theorem 8.42, we know angles θ_1 and θ_2 exist, such that $\mathbf{s} = e^{-i\theta_1}\mathbf{u} > \mathbf{0}$ and $\mathbf{t} = e^{-i\theta_2}\mathbf{v} > \mathbf{0}$. Define $\mathbf{w} = \mathbf{t} - d\mathbf{s}$, where

$$d = \min_{1 \leq j \leq m} s_j^{-1} t_j,$$

so that \mathbf{w} is nonnegative with at least one component equal to 0. If $\mathbf{w} \neq \mathbf{0}$, then clearly $A\mathbf{w} > \mathbf{0}$ because A is positive, which leads to a contradiction because

$$A\mathbf{w} = A\mathbf{t} - dA\mathbf{s} = \rho(A)\mathbf{t} - \rho(A)d\mathbf{s} = \rho(A)\mathbf{w}$$

then implies that $\mathbf{w} > \mathbf{0}$. Thus, we must have $\mathbf{w} = \mathbf{0}$, so $\mathbf{t} = d\mathbf{s}$ and $\mathbf{v} = c\mathbf{u}$, where $c = de^{i(\theta_2 - \theta_1)}$. To prove the second statement of the theorem, first note that from the definition of the spectral radius, $|\lambda| \leq \rho(A)$ for any eigenvalue λ of A . Now if \mathbf{x} is an eigenvector corresponding to λ and $|\lambda| = \rho(A)$, then it follows from Theorem 8.42 that an angle θ exists, such that $\mathbf{u} = e^{-i\theta}\mathbf{x} > \mathbf{0}$. Clearly, $A\mathbf{u} = \lambda\mathbf{u}$. Premultiplying this identity by $D_{\mathbf{u}}^{-1}$, we get

$$D_{\mathbf{u}}^{-1}A\mathbf{u} = \lambda\mathbf{1}_m,$$

so that

$$u_i^{-1} \sum_{j=1}^m a_{ij} u_j = \lambda$$

holds for each i . Now applying Theorem 8.40, we get $\lambda = \rho(A)$. \square

We will see that the first statement in Theorem 8.43 actually can be replaced by the stronger condition that $\rho(A)$ must be a simple eigenvalue of A . However, first we have the following results, the last of which is a useful limiting result for A .

Theorem 8.44 Suppose that A is an $m \times m$ positive matrix, and \mathbf{x} and \mathbf{y} are positive vectors satisfying $A\mathbf{x} = \rho(A)\mathbf{x}$, $A'\mathbf{y} = \rho(A)\mathbf{y}$, and $\mathbf{x}'\mathbf{y} = 1$. Then the following hold:

- (a) $(A - \rho(A)\mathbf{x}\mathbf{y}')^k = A^k - \rho(A)^k\mathbf{x}\mathbf{y}'$, for $k = 1, 2, \dots$.
- (b) Each nonzero eigenvalue of $A - \rho(A)\mathbf{x}\mathbf{y}'$ is an eigenvalue of A .
- (c) $\rho(A)$ is not an eigenvalue of $A - \rho(A)\mathbf{x}\mathbf{y}'$.
- (d) $\rho(A - \rho(A)\mathbf{x}\mathbf{y}') < \rho(A)$.
- (e) $\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A\}^k = \mathbf{x}\mathbf{y}'$.

Proof. (a) is easily established by induction, because it clearly holds for $k = 1$, and if it holds for $k = j - 1$, then

$$\begin{aligned}
 (A - \rho(A)\mathbf{x}\mathbf{y}')^j &= (A - \rho(A)\mathbf{x}\mathbf{y}')^{j-1}(A - \rho(A)\mathbf{x}\mathbf{y}') \\
 &= (A^{j-1} - \rho(A)^{j-1}\mathbf{x}\mathbf{y}')(A - \rho(A)\mathbf{x}\mathbf{y}') \\
 &= A^j - \rho(A)A^{j-1}\mathbf{x}\mathbf{y}' - \rho(A)^{j-1}\mathbf{x}\mathbf{y}'A + \rho(A)^j\mathbf{x}\mathbf{y}'\mathbf{x}\mathbf{y}' \\
 &= A^j - \rho(A)^j\mathbf{x}\mathbf{y}' - \rho(A)^j\mathbf{x}\mathbf{y}' + \rho(A)^j\mathbf{x}\mathbf{y}' \\
 &= A^j - \rho(A)^j\mathbf{x}\mathbf{y}'.
 \end{aligned}$$

Next, suppose that $\lambda \neq 0$ and \mathbf{u} are an eigenvalue and eigenvector of $(A - \rho(A)\mathbf{x}\mathbf{y}')$, so that

$$(A - \rho(A)\mathbf{x}\mathbf{y}')\mathbf{u} = \lambda\mathbf{u}.$$

Premultiplying this equation by $\mathbf{x}\mathbf{y}'$ and observing that $\mathbf{x}\mathbf{y}'(A - \rho(A)\mathbf{x}\mathbf{y}') = \mathbf{0}$, we see that we must have $\mathbf{x}\mathbf{y}'\mathbf{u} = \mathbf{0}$. Consequently,

$$A\mathbf{u} = (A - \rho(A)\mathbf{x}\mathbf{y}')\mathbf{u} = \lambda\mathbf{u},$$

and so λ is also an eigenvalue of A , as is required for (b). To prove (c), suppose that $\lambda = \rho(A)$ is an eigenvalue of $A - \rho(A)\mathbf{x}\mathbf{y}'$ with \mathbf{u} being a corresponding eigenvector. However, we have just seen that this supposition implies that \mathbf{u} is also an eigenvector of A corresponding to the eigenvalue $\rho(A)$. Thus, from Theorem 8.43, $\mathbf{u} = c\mathbf{x}$ for some scalar c and

$$\begin{aligned}
 \rho(A)\mathbf{u} &= (A - \rho(A)\mathbf{x}\mathbf{y}')\mathbf{u} = (A - \rho(A)\mathbf{x}\mathbf{y}')c\mathbf{x} \\
 &= \rho(A)c\mathbf{x} - \rho(A)c\mathbf{x} = \mathbf{0},
 \end{aligned}$$

which is impossible because $\rho(A) > 0$ and $\mathbf{u} \neq \mathbf{0}$, and so (c) holds. Now (d) follows directly from (b), (c), and Theorem 8.43. Finally, to prove (e), note that by dividing both sides of the equation given in (a) by $\rho(A)^k$ and rearranging, we get

$$\{\rho(A)^{-1}A\}^k = \mathbf{x}\mathbf{y}' + \{\rho(A)^{-1}A - \mathbf{x}\mathbf{y}'\}^k.$$

Take the limit, as $k \rightarrow \infty$, of both sides of this equation and observe that from (d),

$$\rho\{\rho(A)^{-1}A - \mathbf{xy}'\} = \frac{\rho\{A - \rho(A)\mathbf{xy}'\}}{\rho(A)} < 1,$$

and so

$$\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A - \mathbf{xy}'\}^k = (0)$$

follows from Theorem 4.26. \square

Theorem 8.45 Let A be an $m \times m$ positive matrix. Then the eigenvalue $\rho(A)$ is a simple eigenvalue of A .

Proof. Let $A = XTX^*$ be the Schur decomposition of A , so that X is a unitary matrix and T is an upper triangular matrix with the eigenvalues of A as its diagonal elements. Write $T = T_1 + T_2$, where T_1 is diagonal and T_2 is upper triangular with each diagonal element equal to 0. Suppose that we have chosen X so that the diagonal elements of T_1 are ordered as $T_1 = \text{diag}(\rho(A), \dots, \rho(A), \lambda_{r+1}, \dots, \lambda_m)$, where r is the multiplicity of the eigenvalue $\rho(A)$ and $|\lambda_j| < \rho(A)$ for $j = r+1, \dots, m$, because of Theorem 8.43. We need to show that $r = 1$. Note that, for any upper triangular matrix U with i th diagonal element u_{ii} , U^k is also upper triangular with its i th diagonal element given by u_{ii}^k . Using this, we find that

$$\begin{aligned} \lim_{k \rightarrow \infty} \{\rho(A)^{-1}A\}^k &= X \left\{ \lim_{k \rightarrow \infty} \{\rho(A)^{-1}(T_1 + T_2)\}^k \right\} X^* \\ &= X \left\{ \lim_{k \rightarrow \infty} \text{diag} \left(1, \dots, 1, \left\{ \frac{\lambda_{r+1}}{\rho(A)} \right\}^k, \dots, \right. \right. \\ &\quad \left. \left. \left\{ \frac{\lambda_m}{\rho(A)} \right\}^k \right) + T_3 \right\} X^* \\ &= X \{\text{diag}(1, \dots, 1, 0, \dots, 0) + T_3\} X^*, \end{aligned}$$

where this last diagonal matrix has r 1's and T_3 is an upper triangular matrix with each diagonal element equal to 0. Clearly, this limiting matrix has rank at least r . However, from Theorem 8.44(e) we see that the limiting matrix must have rank 1, and this proves the result. \square

To this point, we have concentrated on positive matrices. Our next step is to extend some of the previous results to nonnegative matrices. We will see that many of these results generalize to the class of irreducible nonnegative matrices.

Definition 8.2 An $m \times m$ matrix A , with $m \geq 2$, is called a reducible matrix if some integer r , with $1 \leq r \leq m-1$, and $m \times m$ permutation matrix P exist, such that

$$PAP' = \begin{bmatrix} B & C \\ (0) & D \end{bmatrix},$$

where B is $r \times r$, C is $r \times (m - r)$, and D is $(m - r) \times (m - r)$. If A is not reducible, then it is said to be irreducible.

We will need Theorem 8.46 regarding irreducible nonnegative matrices.

Theorem 8.46 An $m \times m$ nonnegative matrix A is irreducible if and only if $(I_m + A)^{m-1} > (0)$.

Proof. First suppose that A is irreducible. We will show that if \mathbf{x} is an $m \times 1$ nonnegative vector with r positive components, $1 \leq r \leq m - 1$, then $(I_m + A)\mathbf{x}$ has at least $r + 1$ positive components. Repeated use of this result verifies that $(I_m + A)^{m-1} > (0)$ because each column of $I_m + A$ has at least two positive components due to the fact that A is irreducible. Since $A \geq (0)$, $(I_m + A)\mathbf{x} = \mathbf{x} + A\mathbf{x}$ must have at least r positive components. If it has exactly r positive components, then the j th component of $A\mathbf{x}$ must be 0 for every j for which $x_j = 0$. Equivalently, for any permutation matrix P , the j th component of $PA\mathbf{x}$ must be 0 for every j for which the j th component of $P\mathbf{x}$ is 0. If we choose a permutation matrix for which $\mathbf{y} = P\mathbf{x}$ has its $m - r$ 0's in the last $m - r$ positions, then we find that the j th component of $PA\mathbf{x} = PAP'\mathbf{y}$ must be 0 for $j = r + 1, \dots, m$. Since $PAP' \geq (0)$ and the first r components of \mathbf{y} are positive, PAP' would have to be of the form

$$PAP' = \begin{bmatrix} B & C \\ (0) & D \end{bmatrix}.$$

Since this result contradicts the fact that A is irreducible, the number of positive components in the vector $(I_m + A)\mathbf{x}$ must exceed r . Conversely, now suppose that $(I_m + A)^{m-1} > (0)$, so that, clearly, $(I_m + A)^{m-1}$ is irreducible. Now A cannot be reducible because, if for some permutation matrix P ,

$$PAP' = \begin{bmatrix} B & C \\ (0) & D \end{bmatrix},$$

then

$$\begin{aligned} P(I_m + A)^{m-1}P' &= \{P(I_m + A)P'\}^{m-1} = (I_m + PAP')^{m-1} \\ &= \begin{bmatrix} I_r + B & C \\ (0) & I_{m-r} + D \end{bmatrix}^{m-1}, \end{aligned}$$

and the matrix on the right-hand side of this last equation has the upper triangular form given in Definition 8.2. \square

We will generalize the result of Theorem 8.41 by showing that $\rho(A)$ is positive, is an eigenvalue of A , and has a positive eigenvector when A is an irreducible nonnegative matrix. However, first we need Theorem 8.47.

Theorem 8.47 Let A be an $m \times m$ irreducible nonnegative matrix, \mathbf{x} be an $m \times 1$ nonnegative vector, and define the function

$$f(\mathbf{x}) = \min_{x_i \neq 0} x_i^{-1} (A)_{i \cdot} \mathbf{x} = \min_{x_i \neq 0} x_i^{-1} \sum_{j=1}^m a_{ij} x_j.$$

Then an $m \times 1$ nonnegative vector \mathbf{b} exists, such that $\mathbf{b}' \mathbf{1}_m = 1$ and $f(\mathbf{b}) \geq f(\mathbf{x})$ holds for any nonnegative \mathbf{x} .

Proof. Define the set

$$S = \{\mathbf{y} : \mathbf{y} = (I_m + A)^{m-1} \mathbf{x}_*, \mathbf{x}_* \in R^m, \mathbf{x}_* \geq \mathbf{0}, \mathbf{x}_*' \mathbf{1}_m = 1\}.$$

Since S is a closed and bounded set, and f is a continuous function on S due to the fact that $\mathbf{y} > \mathbf{0}$ if $\mathbf{y} \in S$, there exists a $\mathbf{c} \in S$, such that $f(\mathbf{c}) \geq f(\mathbf{y})$ for all $\mathbf{y} \in S$. Define $\mathbf{b} = \mathbf{c}/(\mathbf{c}' \mathbf{1}_m)$, and note that f is unaffected by scale changes, so $f(\mathbf{b}) = f(\mathbf{c})$. Let \mathbf{x} be an arbitrary nonnegative vector, and define $\mathbf{x}_* = \mathbf{x}/(\mathbf{x}' \mathbf{1}_m)$ and $\mathbf{y} = (I_m + A)^{m-1} \mathbf{x}_*$. Now it follows from the definition of f that

$$A\mathbf{x}_* - f(\mathbf{x}_*)\mathbf{x}_* \geq \mathbf{0}.$$

Premultiply this equation by $(I_m + A)^{m-1}$ and use the fact that $(I_m + A)^{m-1} A = A(I_m + A)^{m-1}$ to get

$$A\mathbf{y} - f(\mathbf{x}_*)\mathbf{y} \geq \mathbf{0}.$$

However, $\alpha = f(\mathbf{y})$ is the largest value for which $A\mathbf{y} - \alpha\mathbf{y} \geq \mathbf{0}$ because at least one component of $A\mathbf{y} - f(\mathbf{y})\mathbf{y}$ is 0; that is, for some k , $f(\mathbf{y}) = y_k^{-1}(A)_{k \cdot} \mathbf{y}$ and, consequently, the k th component of $A\mathbf{y} - f(\mathbf{y})\mathbf{y}$ will be 0. Thus, we have shown that $f(\mathbf{y}) \geq f(\mathbf{x}_*) = f(\mathbf{x})$. The result then follows from the fact that $f(\mathbf{y}) \leq f(\mathbf{c}) = f(\mathbf{b})$. \square

Theorem 8.48 Let A be an $m \times m$ irreducible nonnegative matrix. Then A has the positive eigenvalue $\rho(A)$ and associated with it a positive eigenvector \mathbf{x} .

Proof. We first show that $f(\mathbf{b})$ is a positive eigenvalue of A , where $f(\mathbf{b})$ is defined as in Theorem 8.47, and \mathbf{b} is a nonnegative vector satisfying $\mathbf{b}' \mathbf{1}_m = 1$ and maximizing f . Since \mathbf{b} maximizes $f(\mathbf{x})$ over all nonnegative \mathbf{x} , we have

$$\begin{aligned} f(\mathbf{b}) &\geq f(m^{-1} \mathbf{1}_m) = \min_{1 \leq i \leq m} (1/m)^{-1} (A)_{i \cdot} (m^{-1} \mathbf{1}_m) \\ &= \min_{1 \leq i \leq m} \sum_{j=1}^m a_{ij} > 0, \end{aligned}$$

because A is nonnegative and irreducible. To prove that $f(\mathbf{b})$ is an eigenvalue of A , recall that from the definition of f it follows that $A\mathbf{b} - f(\mathbf{b})\mathbf{b} \geq \mathbf{0}$. If $A\mathbf{b} - f(\mathbf{b})\mathbf{b}$ has at least one positive component, then because $(I_m + A)^{m-1} > (0)$, we must have

$$(I_m + A)^{m-1}(A\mathbf{b} - f(\mathbf{b})\mathbf{b}) = A\mathbf{y} - f(\mathbf{b})\mathbf{y} > \mathbf{0},$$

where $\mathbf{y} = (I_m + A)^{m-1}\mathbf{b}$. However, $\alpha = f(\mathbf{y})$ is the largest value for which $A\mathbf{y} - \alpha\mathbf{y} \geq \mathbf{0}$, so we would have $f(\mathbf{y}) > f(\mathbf{b})$, which cannot be true because \mathbf{b} maximizes $f(\mathbf{y})$ over all $\mathbf{y} \geq \mathbf{0}$. Thus, $A\mathbf{b} - f(\mathbf{b})\mathbf{b} = \mathbf{0}$, and so $f(\mathbf{b})$ is an eigenvalue of A and \mathbf{b} is a corresponding eigenvector. Our next step is to show that $f(\mathbf{b}) = \rho(A)$ by showing that $f(\mathbf{b}) \geq |\lambda_i|$, where λ_i is an arbitrary eigenvalue of A . Now if \mathbf{u} is an eigenvector of A corresponding to λ_i , then $A\mathbf{u} = \lambda_i\mathbf{u}$ or

$$\lambda_i u_h = \sum_{j=1}^m a_{hj} u_j$$

for $h = 1, \dots, m$. Consequently,

$$|\lambda_i| |u_h| \leq \sum_{j=1}^m a_{hj} |u_j|,$$

for $h = 1, \dots, m$, or simply

$$A \text{ abs}(\mathbf{u}) - |\lambda_i| \text{ abs}(\mathbf{u}) \geq \mathbf{0},$$

which implies that $|\lambda_i| \leq f(\text{abs}(\mathbf{u})) \leq f(\mathbf{b})$. Finally, we must find a positive eigenvector associated with the eigenvalue $\rho(A) = f(\mathbf{b})$. We have already found a nonnegative eigenvector, \mathbf{b} . Note that $A\mathbf{b} = f(\mathbf{b})\mathbf{b}$ implies that $(I_m + A)^{m-1}\mathbf{b} = \{1 + f(\mathbf{b})\}^{m-1}\mathbf{b}$, and so

$$\mathbf{b} = \frac{(I_m + A)^{m-1}\mathbf{b}}{\{1 + f(\mathbf{b})\}^{m-1}}.$$

Thus, using Theorem 8.46, we find that \mathbf{b} is actually positive. □

The proof of Theorem 8.49 will be left to the reader as an exercise.

Theorem 8.49 If A is an $m \times m$ irreducible nonnegative matrix, then $\rho(A)$ is a simple eigenvalue of A .

Although $\rho(A)$ is a simple eigenvalue of an irreducible nonnegative matrix A , there may be other eigenvalues of A that have absolute value $\rho(A)$. Consequently, Theorem 8.44(e) does not immediately extend to irreducible nonnegative matrices. This leads us to the following definition.

Definition 8.3 An $m \times m$ nonnegative matrix A is said to be primitive if it is irreducible and has only one eigenvalue satisfying $|\lambda_i| = \rho(A)$.

Clearly, the result of Theorem 8.44(e) does extend to primitive matrices, and this is summarized in Theorem 8.50.

Theorem 8.50 Let A be an $m \times m$ primitive nonnegative matrix, and suppose that the $m \times 1$ vectors \mathbf{x} and \mathbf{y} satisfy $A\mathbf{x} = \rho(A)\mathbf{x}$, $A'\mathbf{y} = \rho(A)\mathbf{y}$, $\mathbf{x} > \mathbf{0}$, $\mathbf{y} > \mathbf{0}$, and $\mathbf{x}'\mathbf{y} = 1$. Then

$$\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A\}^k = \mathbf{x}\mathbf{y}'.$$

Our final theorem of this section gives a general limit result that holds for all irreducible nonnegative matrices. A proof of this result can be found in Horn and Johnson (2013).

Theorem 8.51 Let A be an $m \times m$ irreducible nonnegative matrix, and suppose that the $m \times 1$ vectors \mathbf{x} and \mathbf{y} satisfy $A\mathbf{x} = \rho(A)\mathbf{x}$, $A'\mathbf{y} = \rho(A)\mathbf{y}$, and $\mathbf{x}'\mathbf{y} = 1$. Then

$$\lim_{N \rightarrow \infty} \left(N^{-1} \sum_{k=1}^N \{\rho(A)^{-1}A\}^k \right) = \mathbf{x}\mathbf{y}'.$$

Nonnegative matrices play an important role in the study of stochastic processes. We will illustrate some of their applications to a particular type of stochastic process known as a Markov chain. Additional information on Markov chains, and stochastic processes in general, can be found in texts such as Bhattacharya and Waymire (2009), Medhi (2009), and Pinsky and Karlin (2011).

Example 8.12 Suppose that we are observing some random phenomenon over time, and at any one point in time our observation can take on any one of the m values, sometimes referred to as states, $1, \dots, m$. In other words, we have a sequence of random variables X_t , for time periods $t = 0, 1, \dots$, where each random variable can be equal to any one of the numbers, $1, \dots, m$. If the probability that X_t is in state i depends only on the state that X_{t-1} is in and not on the states of prior time periods, then this process is said to be a Markov chain. If this probability also does not depend on the value of t , then the Markov chain is said to be homogeneous. In this case, the state probabilities for any time period can be computed from the initial state probabilities and what are known as the transition probabilities. We will write the initial state probability vector $\mathbf{p}^{(0)} = (p_1^{(0)}, \dots, p_m^{(0)})'$, where $p_i^{(0)}$ gives the probability that the process starts out at time 0 in state i . The matrix of transition probabilities is the $m \times m$ matrix P whose (i, j) th element, p_{ij} , gives the probability of X_t being in

state i given that X_{t-1} is in state j . Thus, if $\mathbf{p}^{(t)} = (p_1^{(t)}, \dots, p_m^{(t)})'$ and $p_i^{(t)}$ is the probability that the system is in state i at time t , then, clearly,

$$\begin{aligned}\mathbf{p}^{(1)} &= P\mathbf{p}^{(0)}, \\ \mathbf{p}^{(2)} &= P\mathbf{p}^{(1)} = PP\mathbf{p}^{(0)} = P^2\mathbf{p}^{(0)},\end{aligned}$$

or for general t ,

$$\mathbf{p}^{(t)} = P^t\mathbf{p}^{(0)}.$$

If we have a large population of individuals subject to this random process, then $p_i^{(t)}$ could be described as the proportion of individuals in state i at time t , whereas $p_i^{(0)}$ would be the proportion of individuals starting out in state i . A natural question then is what is happening to these proportions as t increases? That is, can we determine the limiting behavior of $\mathbf{p}^{(t)}$? Note that the answer depends on the limiting behavior of P^t , and P is a nonnegative matrix because each of its elements is a probability. Thus, if P is a primitive matrix, we can apply Theorem 8.50. Now, because the j th column of P gives the probabilities of the various states for time period t when we are in state j at time period $t-1$, the column sum must be 1; that is, $\mathbf{1}'_m P = \mathbf{1}'_m$ or $P'\mathbf{1}_m = \mathbf{1}_m$, so P has an eigenvalue equal to 1. Further, a simple application of Theorem 8.40 assures us that $\rho(P) \leq 1$, so we must have $\rho(P) = 1$. Consequently, if P is primitive and $\boldsymbol{\pi}$ is the $m \times 1$ positive vector satisfying $P\boldsymbol{\pi} = \boldsymbol{\pi}$ and $\boldsymbol{\pi}'\mathbf{1}_m = 1$, then

$$\lim_{t \rightarrow \infty} \{\rho(P)^{-1}P\}^t = \lim_{t \rightarrow \infty} P^t = \boldsymbol{\pi}\mathbf{1}'_m.$$

Using this, we see that

$$\lim_{t \rightarrow \infty} \mathbf{p}^{(t)} = \lim_{t \rightarrow \infty} P^t \mathbf{p}^{(0)} = \boldsymbol{\pi}\mathbf{1}'_m \mathbf{p}^{(0)} = \boldsymbol{\pi},$$

where the last step follows from the fact that $\mathbf{1}'_m \mathbf{p}^{(0)} = 1$. Thus, the system approaches a point of equilibrium in which the proportions for the various states are given by the components of $\boldsymbol{\pi}$, and these proportions do not change from time period to time period. Further, this limiting behavior is not dependent on the initial proportions in $\mathbf{p}^{(0)}$.

As a specific example, let us consider the problem of social mobility that involves the transition between social classes over successive generations in a family. Suppose that each individual is classified according to occupation, as being upper, middle, or lower class, which have been labeled as states 1, 2, and 3, respectively. Suppose that the transition matrix relating a son's class to his father's class is given by

$$P = \begin{bmatrix} 0.45 & 0.05 & 0.05 \\ 0.45 & 0.70 & 0.50 \\ 0.10 & 0.25 & 0.45 \end{bmatrix},$$

so that, for instance, the probabilities that a son will have an upper, middle, and lower class occupation when his father has an upper class occupation are given by the entries in the first column of P . Since P is positive, the limiting result previously discussed applies. A simple eigenanalysis of the matrix P reveals that the positive vector π , which satisfies $P\pi = \pi$ and $\pi' \mathbf{1}_m = 1$, is given by $\pi = (0.083, 0.620, 0.297)'$. Thus, if this random process satisfies the conditions of a homogeneous Markov chain, then after many generations, the male population would consist of 8.3% in the upper class, 62% in the middle class, and 29.7% in the lower class.

8.9 CIRCULANT AND TOEPLITZ MATRICES

In this section, we briefly discuss some structured matrices that have applications in stochastic processes and time series analysis. For a more comprehensive treatment of the first of these classes of matrices, see Davis (1994).

An $m \times m$ matrix A is said to be a circulant matrix if each row of A can be obtained from the previous row by a circular rotation of elements; that is, if we shift each element in the i th row over one column, with the element in the last column being shifted back to the first column, we get the $(i + 1)$ th row, unless $i = m$, in which case we get the first row. Thus, if the elements of the first row of A are a_1, a_2, \dots, a_m , then to be a circulant matrix, A must have the form

$$A = \begin{bmatrix} a_1 & a_2 & a_3 & \cdots & a_{m-1} & a_m \\ a_m & a_1 & a_2 & \cdots & a_{m-2} & a_{m-1} \\ a_{m-1} & a_m & a_1 & \cdots & a_{m-3} & a_{m-2} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_3 & a_4 & a_5 & \cdots & a_1 & a_2 \\ a_2 & a_3 & a_4 & \cdots & a_m & a_1 \end{bmatrix}. \quad (8.20)$$

We will sometimes use the notation $A = \text{circ}(a_1, \dots, a_m)$ to refer to the circulant matrix in (8.20). One special circulant matrix, which we will denote by Π_m , is $\text{circ}(0, 1, 0, \dots, 0)$. This matrix, which also can be written as

$$\Pi_m = (e_m, e_1, \dots, e_{m-1}) = \begin{bmatrix} e'_2 \\ e'_3 \\ \vdots \\ e'_m \\ e'_1 \end{bmatrix},$$

is a permutation matrix, so $\Pi_m^{-1} = \Pi'_m$. Note that if we use $\mathbf{a}_1, \dots, \mathbf{a}_m$ to denote the columns of an arbitrary $m \times m$ matrix A and $\mathbf{b}'_1, \dots, \mathbf{b}'_m$ to denote the rows, then

$$\begin{aligned} A\Pi_m &= (\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_m)(e_m, e_1, \dots, e_{m-1}) \\ &= (\mathbf{a}_m, \mathbf{a}_1, \dots, \mathbf{a}_{m-1}), \end{aligned} \quad (8.21)$$

$$\Pi_m A = \begin{bmatrix} e'_2 \\ e'_3 \\ \vdots \\ e'_m \\ e'_1 \end{bmatrix} \begin{bmatrix} b'_1 \\ b'_2 \\ \vdots \\ b'_{m-1} \\ b'_m \end{bmatrix} = \begin{bmatrix} b'_2 \\ b'_3 \\ \vdots \\ b'_m \\ b'_1 \end{bmatrix}, \quad (8.22)$$

and (8.21) equals (8.22) if and only if A is of the form given in (8.20). Thus, we have Theorem 8.52.

Theorem 8.52 The $m \times m$ matrix A is a circulant matrix if and only if

$$A = \Pi_m A \Pi_m'.$$

Theorem 8.53 gives an expression for an $m \times m$ circulant matrix in terms of a sum of m matrices.

Theorem 8.53 The circulant matrix $A = \text{circ}(a_1, \dots, a_m)$ can be expressed as

$$A = a_1 I_m + a_2 \Pi_m + a_3 \Pi_m^2 + \dots + a_m \Pi_m^{m-1}.$$

Proof. Using (8.20), we see that

$$\begin{aligned} A = & a_1 I_m + a_2 (e_m, e_1, \dots, e_{m-1}) + a_3 (e_{m-1}, e_m, e_1, \dots, e_{m-2}) + \dots \\ & + a_m (e_2, e_3, \dots, e_m, e_1). \end{aligned}$$

Since the postmultiplication of any $m \times m$ matrix by Π_m shifts the columns of that matrix one place to the right, we find that

$$\begin{aligned} \Pi_m^2 &= (e_{m-1}, e_m, \dots, e_{m-2}) \\ &\vdots \\ \Pi_m^{m-1} &= (e_2, e_3, \dots, e_m, e_1), \end{aligned}$$

and so the result follows. \square

Certain operations on circulant matrices produce another circulant matrix. Some of these are given in Theorem 8.54.

Theorem 8.54 Let A and B be $m \times m$ circulant matrices. Then

(a) A' is circulant,

- (b) for any scalars α and β , $\alpha A + \beta B$ is circulant,
- (c) for any positive integer r , A^r is circulant,
- (d) A^{-1} is circulant, if A is nonsingular,
- (e) AB is circulant.

Proof. If $A = \text{circ}(a_1, \dots, a_m)$ and $B = \text{circ}(b_1, \dots, b_m)$, it follows directly from (8.20) that $A' = \text{circ}(a_1, a_m, a_{m-1}, \dots, a_2)$ and

$$\alpha A + \beta B = \text{circ}(\alpha a_1 + \beta b_1, \dots, \alpha a_m + \beta b_m).$$

Since A is circulant, we must have $A = \Pi_m A \Pi_m'$. However, Π_m is an orthogonal matrix, so

$$A^r = (\Pi_m A \Pi_m')^r = \Pi_m A^r \Pi_m',$$

and consequently by Theorem 8.52, A^r is also a circulant matrix. In a similar fashion, we find that if A is nonsingular, then

$$A^{-1} = (\Pi_m A \Pi_m')^{-1} = \Pi_m'^{-1} A^{-1} \Pi_m^{-1} = \Pi_m A^{-1} \Pi_m',$$

and so A^{-1} is circulant. Finally, to prove (e), note that we must have both $A = \Pi_m A \Pi_m'$ and $B = \Pi_m B \Pi_m'$, implying that

$$AB = (\Pi_m A \Pi_m')(\Pi_m B \Pi_m') = \Pi_m AB \Pi_m',$$

and so the proof is complete. \square

The representation of a circulant matrix given in Theorem 8.53 provides a simple way of proving Theorem 8.55.

Theorem 8.55 Suppose that A and B are $m \times m$ circulant matrices. Then their product commutes; that is, $AB = BA$.

Proof. If $A = \text{circ}(a_1, \dots, a_m)$ and $B = \text{circ}(b_1, \dots, b_m)$, then it follows from Theorem 8.53 that

$$A = \sum_{i=1}^m a_i \Pi_m^{i-1}, \quad B = \sum_{j=1}^m b_j \Pi_m^{j-1},$$

where $\Pi_m^0 = I_m$. Consequently,

$$\begin{aligned} AB &= \left(\sum_{i=1}^m a_i \Pi_m^{i-1} \right) \left(\sum_{j=1}^m b_j \Pi_m^{j-1} \right) = \sum_{i=1}^m \sum_{j=1}^m (a_i \Pi_m^{i-1})(b_j \Pi_m^{j-1}) \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i b_j \Pi_m^{i+j-2} = \sum_{i=1}^m \sum_{j=1}^m (b_j \Pi_m^{j-1})(a_i \Pi_m^{i-1}) \end{aligned}$$

$$= \left(\sum_{j=1}^m b_j \Pi_m^{j-1} \right) \left(\sum_{i=1}^m a_i \Pi_m^{i-1} \right) = BA,$$

and so the result follows. \square

All circulant matrices are diagonalizable. We will show this by determining the eigenvalues and eigenvectors of a circulant matrix. However, first let us find the eigenvalues and eigenvectors of the special circulant matrix Π_m .

Theorem 8.56 Let $\lambda_1, \dots, \lambda_m$ be the m solutions to the polynomial equation $\lambda^m - 1 = 0$; that is, $\lambda_j = \theta^{j-1}$, where $\theta = \exp(2\pi i/m) = \cos(2\pi/m) + i \sin(2\pi/m)$ and $i = \sqrt{-1}$. Define Λ to be the diagonal matrix $\text{diag}(1, \theta, \dots, \theta^{m-1})$, and let

$$F = \frac{1}{\sqrt{m}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \theta & \theta^2 & \cdots & \theta^{m-1} \\ 1 & \theta^2 & \theta^4 & \cdots & \theta^{2(m-1)} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & \theta^{m-1} & \theta^{2(m-1)} & \cdots & \theta^{(m-1)(m-1)} \end{bmatrix}.$$

Then the diagonalization of Π_m is given by $\Pi_m = F \Lambda F^*$, where F^* is the conjugate transpose of F ; that is, the diagonal elements of Λ are the eigenvalues of Π_m , whereas the columns of F are corresponding eigenvectors.

Proof. The eigenvalue-eigenvector equation, $\Pi_m \mathbf{x} = \lambda \mathbf{x}$, yields the equations

$$x_{j+1} = \lambda x_j,$$

for $j = 1, \dots, m-1$, and

$$x_1 = \lambda x_m.$$

After repeated substitution, we obtain for any j , $x_j = \lambda^m x_j$. Thus, $\lambda^m = 1$, and so the eigenvalues of Π_m are $1, \theta, \dots, \theta^{m-1}$. Substituting the eigenvalue θ^{j-1} and $x_1 = m^{-1/2}$ into the equations above, we find that an eigenvector corresponding to the eigenvalue θ^{j-1} is given by $\mathbf{x} = m^{-1/2}(1, \theta^{j-1}, \dots, \theta^{(m-1)(j-1)})'$. Thus, we have shown that the diagonal elements of Λ are the eigenvalues of Π_m and the columns of F are corresponding eigenvectors. The remainder of the proof, which simply involves the verification that $F^{-1} = F^*$, is left to the reader as an exercise. \square

The matrix F given in Theorem 8.56 is sometimes referred to as the Fourier matrix of order m . The diagonalization of an arbitrary circulant matrix, which follows directly from Theorem 8.53 and Theorem 8.56, is given in Theorem 8.57.

Theorem 8.57 Let A be the $m \times m$ circulant matrix $\text{circ}(a_1, \dots, a_m)$. Then

$$A = F\Delta F^*,$$

where $\Delta = \text{diag}(\delta_1, \dots, \delta_m)$, $\delta_j = a_1 + a_2\lambda_j^1 + \dots + a_m\lambda_j^{m-1}$ and λ_j and F are defined as in Theorem 8.56.

Proof. Since $\Pi_m = F\Lambda F^*$ and $FF^* = I_m$, we have $\Pi_m^j = F\Lambda^j F^*$, for $j = 2, \dots, m-1$. By using Theorem 8.53, we find that

$$\begin{aligned} A &= a_1 I_m + a_2 \Pi_m + a_3 \Pi_m^2 + \dots + a_m \Pi_m^{m-1} \\ &= a_1 FF^* + a_2 F\Lambda^1 F^* + a_3 F\Lambda^2 F^* + \dots + a_m F\Lambda^{m-1} F^* \\ &= F(a_1 I_m + a_2 \Lambda^1 + a_3 \Lambda^2 + \dots + a_m \Lambda^{m-1})F^* \\ &= F\Delta F^*, \end{aligned}$$

and so the proof is complete. \square

The class of circulant matrices is a subclass of a larger class of matrices known as Toeplitz matrices. The elements of an $m \times m$ Toeplitz matrix A satisfy $a_{ij} = a_{j-i}$ for scalars $a_{-m+1}, a_{-m+2}, \dots, a_{m-1}$; that is, A has the form

$$A = \begin{bmatrix} a_0 & a_1 & a_2 & \cdots & a_{m-2} & a_{m-1} \\ a_{-1} & a_0 & a_1 & \cdots & a_{m-3} & a_{m-2} \\ a_{-2} & a_{-1} & a_0 & \cdots & a_{m-4} & a_{m-3} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_{-m+2} & a_{-m+3} & a_{-m+4} & \cdots & a_0 & a_1 \\ a_{-m+1} & a_{-m+2} & a_{-m+3} & \cdots & a_{-1} & a_0 \end{bmatrix}. \quad (8.23)$$

Two simple $m \times m$ Toeplitz matrices are

$$B = \begin{bmatrix} 0 & 1 & 0 & 0 & \cdots & 0 \\ 0 & 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & 0 & \cdots & 0 \end{bmatrix}, \quad F = \begin{bmatrix} 0 & 0 & \cdots & 0 & 0 & 0 \\ 1 & 0 & \cdots & 0 & 0 & 0 \\ 0 & 1 & \cdots & 0 & 0 & 0 \\ \vdots & \vdots & & \vdots & \vdots & \vdots \\ 0 & 0 & \cdots & 1 & 0 & 0 \\ 0 & 0 & \cdots & 0 & 1 & 0 \end{bmatrix}.$$

These are sometimes referred to as the backward shift and forward shift matrices since for an $m \times m$ matrix $C = [c_1, \dots, c_m]$, $CB = [0, c_1, \dots, c_{m-1}]$ and $CF = [c_2, \dots, c_m, 0]$. Any Toeplitz matrix can be expressed in terms of B and F since $a_0 I_m + \sum_{i=1}^{m-1} (a_{-i} F^i + a_i B^i)$ clearly yields the matrix given in (8.23).

If $a_j = a_{-j}$ for $j = 1, \dots, m-1$, then the matrix A in (8.23) is a symmetric Toeplitz matrix. One important and simple symmetric Toeplitz matrix is one that has $a_j = a_{-j} = 0$ for $j = 2, \dots, m-1$, so that

$$A = \begin{bmatrix} a_0 & a_1 & 0 & \cdots & 0 & 0 \\ a_1 & a_0 & a_1 & \cdots & 0 & 0 \\ 0 & a_1 & a_0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_0 & a_1 \\ 0 & 0 & 0 & \cdots & a_1 & a_0 \end{bmatrix}. \quad (8.24)$$

Toeplitz matrices are sometimes encountered in time series analysis.

Example 8.13 In this example, we consider a response variable y which is observed over time. For instance, we may be recording the annual amount of rainfall at some particular location and have the measurements for the last N years. These type of data are called time series, and standard regression methods typically cannot be used since the responses are usually correlated. One of the time series models sometimes used for this type of data is a moving average model which models the response variable for time period i as

$$y_i = \sum_{j=0}^{\infty} \rho^j \epsilon_{i-j}.$$

Here ρ is a constant satisfying $|\rho| < 1$, while the ϵ_{i-j} 's are uncorrelated random errors each with mean 0 and variance σ^2 . We will determine the covariance matrix of the vector, $\mathbf{y} = (y_1, \dots, y_N)'$, of the N successive observations of the response. Clearly, $E(y_i) = 0$ for all i so for $i = 1, \dots, N$ and $h = 0, \dots, N-i$, we have

$$\begin{aligned} \text{cov}(y_i, y_{i+h}) &= E(y_i y_{i+h}) = E \left(\sum_{j=0}^{\infty} \sum_{l=0}^{\infty} \rho^j \epsilon_{i-j} \rho^l \epsilon_{i+h-l} \right) \\ &= E \left(\sum_{j=0}^{\infty} \sum_{k=-h}^{\infty} \rho^{j+k+h} \epsilon_{i-j} \epsilon_{i-k} \right) = \sum_{j=0}^{\infty} \rho^{2j+h} \sigma^2 \\ &= \rho^h \sigma^2 \sum_{j=0}^{\infty} \rho^{2j} = \frac{\rho^h \sigma^2}{1 - \rho^2} = \rho^h \gamma^2, \end{aligned}$$

where $\gamma^2 = \text{var}(y_i) = \sigma^2/(1 - \rho^2)$. Thus, the covariance matrix of \mathbf{y} is the symmetric Toeplitz matrix

$$\gamma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \cdots & \rho^{N-1} \\ \rho & 1 & \rho & \cdots & \rho^{N-2} \\ \vdots & \vdots & \vdots & & \vdots \\ \rho^{N-1} & \rho^{N-2} & \rho^{N-3} & \cdots & 1 \end{bmatrix}.$$

Some specialized results, such as formulas for eigenvalues and formulas for the computation of the inverse of a Toeplitz matrix, can be found in Grenander and Szego (1984) and Heinig and Rost (1984).

8.10 HADAMARD AND VANDERMONDE MATRICES

In this section, we discuss some matrices that have applications in the areas of design of experiments and response surface methodology. We begin with a class of matrices known as Hadamard matrices. An $m \times m$ matrix H is said to be a Hadamard matrix if first, each element of H is either $+1$ or -1 , and second, H satisfies

$$H'H = HH' = mI_m; \quad (8.25)$$

that is, the columns of H form an orthogonal set of vectors, and the rows form an orthogonal set as well. For instance, a 2×2 Hadamard matrix is given by

$$H = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix},$$

whereas a 4×4 Hadamard matrix is given by

$$H = \begin{bmatrix} 1 & 1 & 1 & 1 \\ -1 & -1 & 1 & 1 \\ 1 & -1 & 1 & -1 \\ 1 & -1 & -1 & 1 \end{bmatrix}.$$

Some of the basic properties of Hadamard matrices are given in Theorem 8.58.

Theorem 8.58 Let H_m denote any $m \times m$ Hadamard matrix. Then

- (a) $m^{-1/2}H_m$ is an $m \times m$ orthogonal matrix,
- (b) $|H_m| = \pm m^{m/2}$,
- (c) $H_m \otimes H_n$ is an $mn \times mn$ Hadamard matrix.

Proof. (a) follows directly from (8.25). Also using (8.25), we find that

$$|H'_m H_m| = |mI_m| = m^m.$$

However,

$$|H'_m H_m| = |H'_m| |H_m| = |H_m|^2,$$

and so (b) follows. Note that each element of $H_m \otimes H_n$ is $+1$ or -1 because each element is the product of an element from H_m and an element from H_n , and

$$\begin{aligned}(H_m \otimes H_n)'(H_m \otimes H_n) &= H_m' H_m \otimes H_n' H_n \\ &= mI_m \otimes nI_n = mnI_{mn},\end{aligned}$$

so (c) follows. \square

Hadamard matrices that have all of the elements of the first row equal to $+1$ are called normalized Hadamard matrices. Our next result addresses the existence of normalized Hadamard matrices.

Theorem 8.59 If an $m \times m$ Hadamard matrix exists, then an $m \times m$ normalized Hadamard matrix exists.

Proof. Suppose that H is an $m \times m$ Hadamard matrix. Let D be the diagonal matrix with the elements of the first row of H as its diagonal elements; that is, $D = \text{diag}(h_{11}, \dots, h_{1m})$. Note that $D^2 = I_m$ because each diagonal element of D is $+1$ or -1 . Consider the $m \times m$ matrix $H_* = HD$. Each column of H_* is the corresponding column of H multiplied by either $+1$ or -1 , so clearly each element of H_* is $+1$ or -1 . The j th element in the first row of H_* is $h_{1j}^2 = 1$, so H_* has all of its elements of the first row equal to $+1$. In addition,

$$\begin{aligned}H_*' H_* &= (HD)' HD = D' H' HD \\ &= D(mI_m)D = mD^2 = mI_m.\end{aligned}$$

Thus, H_* is an $m \times m$ normalized Hadamard matrix, and so the proof is complete. \square

Hadamard matrices of size $m \times m$ do not exist for every choice of m . We have already given an example of a 2×2 Hadamard matrix, and this matrix can be used repeatedly in Theorem 8.58(c) to obtain a $2^n \times 2^n$ Hadamard matrix for any integer $n \geq 2$. However, $m \times m$ Hadamard matrices do exist for some values of $m \neq 2^n$. Theorem 8.60 gives a necessary condition on the order m so that Hadamard matrices of order m exist.

Theorem 8.60 If H is an $m \times m$ Hadamard matrix, where $m > 2$, then m is a multiple of four.

Proof. The result can be proven by using the fact that any three rows of H are orthogonal to one another. Consequently, we will refer to the first three rows of H , and, because of Theorem 8.59, we may assume that H is a normalized Hadamard matrix, so that all of the elements in the first row are $+1$. Since the second and third

rows are orthogonal to the first row, they must each have $r + 1$'s and $r - 1$'s, where $r = m/2$; thus, clearly,

$$m = 2r, \quad (8.26)$$

or, in other words, m is a multiple of 2. Let n_{+-} be the number of columns in which row 2 has a $+1$ and row 3 has a -1 . Similarly, define n_{-+} , n_{++} , and n_{--} . Note that the value of any one of these n 's determines the others because $n_{++} + n_{+-} = r$, $n_{++} + n_{-+} = r$, $n_{--} + n_{-+} = r$, and $n_{--} + n_{+-} = r$. For instance, if $n_{++} = s$, then $n_{+-} = (r - s)$, $n_{-+} = (r - s)$, and $n_{--} = s$. However, the orthogonality of rows 2 and 3 guarantee that $n_{++} + n_{--} = n_{-+} + n_{+-}$, which yields the relationship

$$2s = 2(r - s).$$

Thus, $r = 2s$, and so using (8.26) we get $m = 4s$, which completes the proof. \square

Some additional results on Hadamard matrices can be found in Hedayat and Wallis (1978), Agaian (1985), and Xian (2001).

An $m \times n$ matrix A is said to be a Vandermonde matrix if it has the form

$$A = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ a_1 & a_2 & a_3 & \cdots & a_n \\ a_1^2 & a_2^2 & a_3^2 & \cdots & a_n^2 \\ \vdots & \vdots & \vdots & & \vdots \\ a_1^{m-1} & a_2^{m-1} & a_3^{m-1} & \cdots & a_n^{m-1} \end{bmatrix}. \quad (8.27)$$

For instance, if F is the $m \times m$ Fourier matrix discussed in Section 8.9, then $A = m^{1/2}F$ is an $m \times m$ Vandermonde matrix with $a_i = \theta^{i-1}$, for $i = 1, \dots, m$. For a statistical example, consider the polynomial regression model,

$$y_i = \beta_0 + \beta_1 x_i + \cdots + \beta_k x_i^k + \epsilon_i,$$

in which the response variable y is regressed on one explanatory variable x through a k th-degree polynomial. If we have N observations and the model is written in the usual matrix form $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, then X' has the form of the Vandermonde matrix in (8.27) with $a_i = x_i$, $m = k + 1$, and $n = N$.

Our final result of this chapter gives an expression for the determinant of a square Vandermonde matrix.

Theorem 8.61 Let A be the $m \times m$ Vandermonde matrix given in (8.27). Then its determinant is given by

$$|A| = \prod_{1 \leq i < j \leq m} (a_j - a_i). \quad (8.28)$$

Proof. Our proof is by induction. For $m = 2$, we find that

$$|A| = \begin{vmatrix} 1 & 1 \\ a_1 & a_2 \end{vmatrix} = a_2 - a_1,$$

and so (8.28) holds when A is 2×2 . Next we assume that (8.28) holds for Vandermonde matrices of order $m - 1$ and show that it must also hold for order m . Thus, if B is the $(m - 1) \times (m - 1)$ matrix obtained from A by deleting its last row and first column, then, because B is a Vandermonde matrix of order $m - 1$, we must have

$$|B| = \prod_{2 \leq i < j \leq m} (a_j - a_i).$$

Define the $m \times m$ matrix

$$C = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 & 0 \\ -a_1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & -a_1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & -a_1 & 1 \end{bmatrix},$$

and note that by repeatedly using the cofactor expansion formula for a determinant on the first row, we find that $|C| = 1$. Thus, $|A| = |CA|$. However, it is easily verified that $CA = E$, where

$$E = \begin{bmatrix} 1 & \mathbf{1}'_{m-1} \\ \mathbf{0} & BD \end{bmatrix},$$

and $D = \text{diag}((a_2 - a_1), (a_3 - a_1), \dots, (a_m - a_1))$. Consequently,

$$\begin{aligned} |A| &= |CA| = |E| = |BD| = |B||D| \\ &= \left\{ \prod_{2 \leq i < j \leq m} (a_j - a_i) \right\} \left\{ \prod_{2 \leq j \leq m} (a_j - a_1) \right\} \\ &= \prod_{1 \leq i < j \leq m} (a_j - a_i), \end{aligned}$$

where the third equality was obtained by using Theorem 7.4. This completes the proof. \square

PROBLEMS

8.1 Let the 2×2 matrices A and B be given by

$$A = \begin{bmatrix} 2 & 3 \\ 1 & 2 \end{bmatrix}, \quad B = \begin{bmatrix} 5 & 3 \\ 3 & 2 \end{bmatrix}.$$

- (a) Compute $A \otimes B$ and $B \otimes A$.
 - (b) Find $\text{tr}(A \otimes B)$.
 - (c) Compute $|A \otimes B|$.
 - (d) Give the eigenvalues of $A \otimes B$.
 - (e) Find $(A \otimes B)^{-1}$.
- 8.2** Give a simplified expression for $I_m \otimes I_n$.
- 8.3** Prove the properties given in Theorem 8.1.
- 8.4** Suppose that A and B are $m \times n$ and $p \times q$ matrices, respectively, and c is an $r \times 1$ vector. Show that
- (a) $A(I_n \otimes c') = A \otimes c'$,
 - (b) $(c \otimes I_p)B = c \otimes B$.
- 8.5** Let a be an $m \times 1$ vector and B be a $p \times q$ matrix. Suppose B is partitioned as $B = [B_1 \cdots B_k]$. Show that

$$a \otimes B = [a \otimes B_1 \cdots a \otimes B_k].$$

- 8.6** Prove results (b) and (c) of Theorem 8.4.
- 8.7** Show that if A and B are symmetric matrices, then $A \otimes B$ is also symmetric.
- 8.8** Show that $A \otimes B$ is nonsingular if and only if A and B are nonsingular.
- 8.9** Let A be $m \times m$ and B be $n \times n$. Show that $A \otimes B$ is an orthogonal matrix if and only if cA and $c^{-1}B$ are orthogonal matrices for some $c > 0$.
- 8.10** Find the rank of $A \otimes B$, where

$$A = \begin{bmatrix} 2 & 6 \\ 1 & 4 \\ 3 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 5 & 2 & 4 \\ 2 & 1 & 1 \\ 1 & 0 & 2 \end{bmatrix}.$$

- 8.11** For matrices A and B of any sizes, show that $A \otimes B = (0)$ if and only if $A = (0)$ or $B = (0)$.

- 8.12** Let \mathbf{x}_i be an eigenvector of the $m \times m$ matrix A corresponding to the eigenvalue λ_i . Let \mathbf{y}_j be an eigenvector of the $p \times p$ matrix B corresponding to the eigenvalue θ_j .
- (a) Show that $\mathbf{x}_i \otimes \mathbf{y}_j$ is an eigenvector of $A \otimes B$.
- (b) Give an example of matrices A and B , such that $A \otimes B$ has an eigenvector that is not the Kronecker product of an eigenvector of A and an eigenvector of B .
- 8.13** Suppose that A is an $m \times n$ matrix and B is an $n \times m$ matrix, where $n > m$. Show that $A \otimes B$ has a 0 eigenvalue with multiplicity at least $(n - m)m$.
- 8.14** Suppose \mathbf{x} is an eigenvector of the $m \times m$ matrix A corresponding to the eigenvalue λ , and \mathbf{y} is an eigenvector of the $n \times n$ matrix B corresponding to the eigenvalue μ . Show that $\mathbf{y} \otimes \mathbf{x}$ is an eigenvector of $(I_n \otimes A) + (B \otimes I_m)$ corresponding to the eigenvalue $\lambda + \mu$.
- 8.15** Show that if A and B are positive definite matrices, then $A \otimes B$ is also positive definite.
- 8.16** It follows from Theorem 8.3 and Theorem 8.6 that if A and B are square matrices, then $\text{tr}(A \otimes B) = \text{tr}(B \otimes A)$ and $|A \otimes B| = |B \otimes A|$. Show that when A and B are not square and $A \otimes B$ is square, then the first of these two identities need not hold whereas the second one does hold. That is, suppose that A is $m \times n$ and B is $n \times m$.
- (a) Give an example for which $\text{tr}(A \otimes B) \neq \text{tr}(B \otimes A)$.
- (b) Prove that $|A \otimes B| = |B \otimes A|$.
- 8.17** Let \mathbf{x} be an $m \times 1$ vector and \mathbf{y} be an $n \times 1$ vector. Verify that the three matrices $\mathbf{x}\mathbf{y}'$, $\mathbf{y}'\mathbf{x}$, and $\mathbf{x} \otimes \mathbf{y}'$ are identical.
- 8.18** Compute the sum of squared errors $\text{SSE} = (\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}})$ for the two-way classification model with interaction discussed in Example 8.3.
- 8.19** Consider the two-way classification model without interaction given by

$$y_{ijk} = \mu + \tau_i + \gamma_j + \epsilon_{ijk},$$

where $i = 1, \dots, a$, $j = 1, \dots, b$, and $k = 1, \dots, n$.

- (a) Find a least squares solution for $\boldsymbol{\beta} = (\mu, \tau_1, \dots, \tau_a, \gamma_1, \dots, \gamma_b)'$, and use this to obtain the vector of fitted values and the sum of squared errors for this model.
- (b) Compute the sum of squared errors for the reduced model $y_{ijk} = \mu + \gamma_j + \epsilon_{ijk}$ and use this along with the SSE computed in (a) to show that the sum of squares for factor A is

$$\text{SSA} = nb \sum_{i=1}^a (\bar{y}_{i.} - \bar{y}_{..})^2.$$

- (c) In a similar fashion, show that the sum of squares for factor B is

$$\text{SSB} = na \sum_{j=1}^b (\bar{y}_{.j} - \bar{y}_{..})^2.$$

- (d) Find a set of as many linearly independent estimable functions of μ , τ_i , and γ_j as possible.
- (e) Use the sum of squared errors computed in (a) and the sum of squared errors computed in Problem 8.18 to show that the sum of squares for interaction in the model of Problem 8.18 is given by

$$\text{SSAB} = n \sum_{i=1}^a \sum_{j=1}^b (\bar{y}_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2.$$

8.20 Prove Theorem 8.8.

8.21 Let A_1, A_2, A_3 , and A_4 be square matrices. Show that, when the sizes of these matrices are such that the appropriate operations are defined,

(a) $(A_1 \oplus A_2) + (A_3 \oplus A_4) = (A_1 + A_3) \oplus (A_2 + A_4),$

(b) $(A_1 \oplus A_2)(A_3 \oplus A_4) = A_1 A_3 \oplus A_2 A_4,$

(c) $(A_1 \oplus A_2) \otimes A_3 = (A_1 \otimes A_3) \oplus (A_2 \otimes A_3).$

8.22 Give an example to show that, in general,

$$A_1 \otimes (A_2 \oplus A_3) \neq (A_1 \otimes A_2) \oplus (A_1 \otimes A_3).$$

8.23 Use Theorem 6.4 and Theorem 8.11 to prove Theorem 6.5.

8.24 Consider the system of equations $AX - XB = C$, where X is an $m \times n$ matrix of variables and A, B , and C are matrices of constants. Show that if the matrices A and B do not have any eigenvalues in common, then this system has a unique solution for X .

8.25 Prove the results of Corollary 8.12.1.

8.26 Let A and B be $m \times n$ and $n \times p$ matrices, respectively, whereas c and d are $p \times 1$ and $n \times 1$ vectors, respectively. Show that

(a) $ABc = (c' \otimes A) \text{vec}(B) = (A \otimes c') \text{vec}(B'),$

(b) $d'Bc = (c' \otimes d') \text{vec}(B).$

8.27 Let A, B , and C be $m \times m$ matrices. Show that if C is symmetric, then

$$\{\text{vec}(C)\}'(A \otimes B) \text{vec}(C) = \{\text{vec}(C)\}'(B \otimes A) \text{vec}(C).$$

8.28 For any matrix A and any vector b , show that

$$\text{vec}(A \otimes b) = \text{vec}(A) \otimes b.$$

8.29 Let A be an $m \times n$ matrix and B be an $n \times p$ matrix. Show that

$$\begin{aligned}\text{vec}(AB) &= (I_p \otimes A) \text{vec}(B) = (B' \otimes I_m) \text{vec}(A) \\ &= (B' \otimes A) \text{vec}(I_n).\end{aligned}$$

8.30 Let A be an $m \times m$ matrix, B be an $n \times n$ matrix, and C be an $m \times n$ matrix. Prove that

$$\text{vec}(AC + CB) = \{(I_n \otimes A) + (B' \otimes I_m)\} \text{vec}(C).$$

8.31 Let A and B be $m \times n$ matrices. Show that

$$\{\text{tr}(A'B)\}^2 \leq \{\text{tr}(A'A)\}\{\text{tr}(B'B)\}$$

with equality if and only if one of the matrices is a scalar multiple of the other.

8.32 Let A be an $m \times m$ symmetric matrix, and consider the function of A defined by $f(A) = \text{tr}(A^2) - m^{-1}\{\text{tr}(A)\}^2$. Show that $f(A)$ can be expressed as

$$f(A) = \{\text{vec}(A)\}' \{I_{m^2} - m^{-1} \text{vec}(I_m) \text{vec}(I_m)'\} \text{vec}(A).$$

8.33 If e_i is the i th column of the identity matrix I_m , verify that

$$\text{vec}(I_m) = \sum_{i=1}^m (e_i \otimes e_i).$$

8.34 Prove property (h) of Theorem 8.13.

8.35 Let the 2×2 matrices A and B be given by

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}, \quad B = \begin{bmatrix} 4 & 1 \\ 1 & 3 \end{bmatrix}.$$

(a) Compute $A \odot B$.

(b) Which of the matrices, A , B , and $A \odot B$, are positive definite or positive semidefinite? How does this relate to Theorem 8.17?

8.36 Give an example of matrices A and B such that neither is nonnegative definite, yet $A \odot B$ is positive definite.

8.37 Let A , B , and C be $m \times n$ matrices. Show that

$$\text{tr}\{(A' \odot B')C\} = \text{tr}\{A'(B \odot C)\}.$$

- 8.38** Suppose that the $m \times m$ matrix A is diagonalizable; that is, a nonsingular matrix X and a diagonal matrix $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ exist, such that $A = X\Lambda X^{-1}$. Show that if we define the vector of diagonal elements of A , $\mathbf{a} = (a_{11}, \dots, a_{mm})'$, and the vector of eigenvalues of A , $\boldsymbol{\lambda} = (\lambda_1, \dots, \lambda_m)'$, then

$$(X \odot X^{-1'})\boldsymbol{\lambda} = \mathbf{a}$$

and

$$(X \odot X^{-1'})\mathbf{1}_m = (X \odot X^{-1'})'\mathbf{1}_m = \mathbf{1}_m.$$

- 8.39** Let A and B be $m \times m$ nonnegative definite matrices. Show that

(a) $|A \odot B| \geq |A||B|$,

(b) $|A \odot A^{-1}| \geq 1$, if A is positive definite.

- 8.40** Let A be an $m \times m$ matrix. Use Theorem 8.15 to show that A is nonnegative definite if and only if $\text{tr}(AB') \geq 0$ for every $m \times m$ nonnegative definite matrix B .

- 8.41** For each of the following pairs of 2×2 matrices, compute the smaller eigenvalue $\lambda_2(A \odot B)$ and the lower bounds for this eigenvalue given by Theorem 8.21 and Theorem 8.23. Which bound is closer to the actual value?

(a) $A = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}.$

(b) $A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & \sqrt{2} \\ \sqrt{2} & 3 \end{bmatrix}.$

- 8.42** Let A be an $m \times m$ positive definite matrix. Use Theorem 8.19 to show that if $B = A^{-1}$, then $a_{11}b_{11} \geq 1$. Show how this result generalizes to $a_{ii}b_{ii} \geq 1$ for $i = 1, \dots, m$.

- 8.43** Let A and B be $m \times m$ positive definite matrices, and consider the inequality

$$|A \odot B| + |A||B| \geq |A| \prod_{i=1}^m b_{ii} + |B| \prod_{i=1}^m a_{ii}.$$

- (a) Show that this inequality is equivalent to

$$|R_A \odot R_B| + |R_A||R_B| \geq |R_A| + |R_B|,$$

where R_A and R_B represent the correlation matrices computed from A and B .

- (b) Use Theorem 8.20 on $|R_A \odot C|$, where $C = R_B - (\mathbf{e}_1' R_B^{-1} \mathbf{e}_1)^{-1} \mathbf{e}_1 \mathbf{e}_1'$, to establish the inequality given in (a).

- 8.44** Suppose that A and B are $m \times m$ positive definite matrices. Show that $A \odot B = AB$ if and only if both A and B are diagonal matrices.
- 8.45** Let A be an $m \times m$ positive definite matrix and B be an $m \times m$ positive semidefinite matrix with exactly r positive diagonal elements. Show that $\text{rank}(A \odot B) = r$.
- 8.46** Show that if A and B are singular 2×2 matrices, then $A \odot B$ is also singular.
- 8.47** Let R be an $m \times m$ positive definite correlation matrix having λ as its smallest eigenvalue. Show that if τ is the smallest eigenvalue of $R \odot R$ and $R \neq I_m$, then $\tau > \lambda$.
- 8.48** Consider the matrix

$$\Psi_m = \sum_{i=1}^m \mathbf{e}_{i,m} (\mathbf{e}_{i,m} \otimes \mathbf{e}_{i,m})',$$

which we have seen satisfies $\Psi_m(A \otimes B)\Psi_m' = A \odot B$ for any $m \times m$ matrices A and B . Define $\mathbf{w}(A)$ to be the $m \times 1$ vector containing the diagonal elements of A ; that is, $\mathbf{w}(A) = (a_{11}, \dots, a_{mm})'$. Also let Λ_m be the $m^2 \times m^2$ matrix given by

$$\Lambda_m = \sum_{i=1}^m (E_{ii} \otimes E_{ii}) = \sum_{i=1}^m (\mathbf{e}_{i,m} \mathbf{e}_{i,m}' \otimes \mathbf{e}_{i,m} \mathbf{e}_{i,m}').$$

Show that

- (a) $\Psi_m' \mathbf{w}(A) = \text{vec}(A)$ for every diagonal matrix A ,
- (b) $\Psi_m \text{vec}(A) = \mathbf{w}(A)$ for every matrix A ,
- (c) $\Psi_m \Psi_m' = I_m$ so that $\Psi_m^+ = \Psi_m'$,
- (d) $\Psi_m' \Psi_m = \Lambda_m$,
- (e) $\Lambda_m N_m = N_m \Lambda_m = \Lambda_m$,
- (f) $\{\text{vec}(A)\}' \Lambda_m (B \otimes B) \Lambda_m \text{vec}(A) = \{\mathbf{w}(A)\}' (B \odot B) \mathbf{w}(A)$.

Additional properties of Ψ_m can be found in Magnus (1988).

- 8.49** Let A and B be $m \times m$ positive definite matrices. Since $\Psi_m(A \otimes B)\Psi_m' = A \odot B$ and $\Psi_m \Psi_m' = I_m$, it follows that an $m^2 \times m^2$ orthogonal matrix P exists, such that $P(A \otimes B)P'$ can be partitioned into the 2×2 form of (7.1) with the $(1, 1)$ th submatrix given by $A \odot B$. Use this result and the result from Problem 7.10 to show that
- (a) $A^{-1} \odot B^{-1} - (A \odot B)^{-1}$ is nonnegative definite,
 - (b) $A^{-1} \odot A^{-1} - (A \odot A)^{-1}$ is nonnegative definite,
 - (c) $A^{-1} \odot A - (A^{-1} \odot A)^{-1}$ is nonnegative definite.
- 8.50** Verify that the commutation matrix K_{mn} is a permutation matrix; that is, show that each column of K_{mn} is a column of I_{mn} and each column of I_{mn} is a column of K_{mn} .
- 8.51** Write out the commutation matrices K_{22} and K_{24} .

- 8.52** The eigenvalues of K_{mn} were given in Theorem 8.28. Show that corresponding eigenvectors are given by the vectors of the form $\mathbf{e}_l \otimes \mathbf{e}_l$, $(\mathbf{e}_l \otimes \mathbf{e}_k) + (\mathbf{e}_k \otimes \mathbf{e}_l)$, and $(\mathbf{e}_l \otimes \mathbf{e}_k) - (\mathbf{e}_k \otimes \mathbf{e}_l)$.
- 8.53** Show that the commutation matrix K_{mn} can be expressed as

$$K_{mn} = \sum_{i=1}^m (\mathbf{e}_i \otimes I_n \otimes \mathbf{e}_i'),$$

where \mathbf{e}_i is the i th column of I_m . Use this to show that if A is $n \times m$, \mathbf{x} is $m \times 1$, and \mathbf{y} is an arbitrary vector, then

$$K'_{mn}(\mathbf{x} \otimes A \otimes \mathbf{y}') = A \otimes \mathbf{x}\mathbf{y}'.$$

- 8.54** Show that

$$(a) \quad K_{np,m} = K_{n,pm} K_{p,nm} = K_{p,nm} K_{n,pm},$$

$$(b) \quad K_{np,m} K_{pm,n} K_{mn,p} = I_{mnp}.$$

- 8.55** Let A be an $m \times n$ matrix with rank r , and let $\lambda_1, \dots, \lambda_r$ be the nonzero eigenvalues of $A'A$. If we define

$$P = K_{mn}(A' \otimes A),$$

show that

$$(a) \quad P \text{ is symmetric,}$$

$$(b) \quad \text{rank}(P) = r^2,$$

$$(c) \quad \text{tr}(P) = \text{tr}(A'A),$$

$$(d) \quad P^2 = (AA') \otimes (A'A),$$

$$(e) \quad \text{the nonzero eigenvalues of } P \text{ are } \lambda_1, \dots, \lambda_r \text{ and } \pm(\lambda_i \lambda_j)^{1/2} \text{ for all } i < j.$$

- 8.56** Let A be an $m \times n$ matrix and B be a $p \times q$ matrix. Show that

$$(a) \quad \text{vec}(A' \otimes B) = (K_{mq,n} \otimes I_p) \{\text{vec}(A) \otimes \text{vec}(B)\},$$

$$(b) \quad \text{vec}(A \otimes B') = (I_n \otimes K_{p,mq}) \{\text{vec}(A) \otimes \text{vec}(B)\}.$$

- 8.57** Suppose that A is an $m \times n$ matrix and B is a $p \times q$ matrix with $mp = nq$. Show that

$$\text{tr}(A \otimes B) = \{\text{vec}(I_n) \otimes \text{vec}(I_q)\}' \{\text{vec}(A) \otimes \text{vec}(B')\}.$$

- 8.58** Show that

$$(a) \quad K_{mnp,q} = (I_{mn} \otimes K_{pq})(I_m \otimes K_{nq} \otimes I_p)(K_{mq} \otimes I_{np}),$$

$$(b) \quad K_{mn,pq} = (I_m \otimes K_{np} \otimes I_q)(K_{mp} \otimes K_{nq})(I_p \otimes K_{mq} \otimes I_n).$$

- 8.59** Prove the results of Theorem 8.31.

8.60 Show that if A and B are $m \times m$ matrices, then

$$\begin{aligned} N_m(A \otimes B + B \otimes A)N_m &= (A \otimes B + B \otimes A)N_m \\ &= N_m(A \otimes B + B \otimes A) \\ &= 2N_m(A \otimes B)N_m. \end{aligned}$$

8.61 Consider the matrix $N_m = \frac{1}{2}(I_{m^2} + K_{mm})$.

(a) Show that N_m can be expressed as

$$N_m = \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (e_i e'_i \otimes e_j e'_j + e_i e'_j \otimes e_j e'_i).$$

(b) Show that

$$N_m(\mathbf{a} \otimes \mathbf{b}) = \frac{1}{2}(\mathbf{a} \otimes \mathbf{b} + \mathbf{b} \otimes \mathbf{a})$$

for any $m \times 1$ vectors \mathbf{a} and \mathbf{b} .

(c) Let Δ be the matrix that generalizes the property illustrated in (b) to the Kronecker product of three $m \times 1$ vectors, $\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}$; that is, suppose Δ satisfies

$$\begin{aligned} \Delta(\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c}) &= \frac{1}{6}(\mathbf{a} \otimes \mathbf{b} \otimes \mathbf{c} + \mathbf{a} \otimes \mathbf{c} \otimes \mathbf{b} + \mathbf{b} \otimes \mathbf{a} \otimes \mathbf{c} \\ &\quad + \mathbf{b} \otimes \mathbf{c} \otimes \mathbf{a} + \mathbf{c} \otimes \mathbf{a} \otimes \mathbf{b} + \mathbf{c} \otimes \mathbf{b} \otimes \mathbf{a}). \end{aligned}$$

Show that Δ can be expressed as

$$\begin{aligned} \Delta &= \frac{1}{6} \sum_{h=1}^m \sum_{i=1}^m \sum_{j=1}^m (e_h e'_h \otimes e_i e'_i \otimes e_j e'_j + e_h e'_h \otimes e_i e'_j \otimes e_j e'_i \\ &\quad + e_h e'_i \otimes e_i e'_h \otimes e_j e'_j + e_h e'_j \otimes e_i e'_h \otimes e_j e'_i \\ &\quad + e_h e'_i \otimes e_i e'_j \otimes e_j e'_h + e_h e'_j \otimes e_i e'_i \otimes e_j e'_h). \end{aligned}$$

8.62 Write out the matrices N_2 and N_3 .

8.63 For $i = 1, \dots, m, j = 1, \dots, i$, define the $m(m+1)/2 \times 1$ vector \mathbf{u}_{ij} to be the vector with one in its $\{(j-1)m + i - j(j-1)/2\}$ th position and zeros elsewhere. It can be easily verified that these vectors are the columns of the identity matrix of order $m(m+1)/2$; that is,

$$I_{m(m+1)/2} = (\mathbf{u}_{11}, \mathbf{u}_{21}, \dots, \mathbf{u}_{m1}, \mathbf{u}_{22}, \dots, \mathbf{u}_{m2}, \mathbf{u}_{33}, \dots, \mathbf{u}_{mm}).$$

Let E_{ij} be the $m \times m$ matrix whose only nonzero element is a one in the (i, j) th position, and define

$$T_{ij} = \begin{cases} E_{ij} + E_{ji}, & \text{if } i \neq j, \\ E_{ii}, & \text{if } i = j. \end{cases}$$

Show that $D_m = \sum_{i \geq j} \{\text{vec}(T_{ij})\} \mathbf{u}'_{ij}$; that is, verify that

$$\sum_{i \geq j} \{\text{vec}(T_{ij})\} \mathbf{u}'_{ij} \mathbf{v}(A) = \text{vec}(A),$$

where A is an arbitrary $m \times m$ symmetric matrix.

8.64 Use the expression given for D_m in Problem 8.63 to prove Theorem 8.34(b).

8.65 Let \mathbf{u}_{ij} , $i = 1, \dots, m$, $j = 1, \dots, i$ be the $m(m+1)/2 \times 1$ vectors defined in Problem 8.63. Show that

$$(a) \quad D'_m D_m = 2I_{m(m+1)/2} - \sum_{i=1}^m \mathbf{u}_{ii} \mathbf{u}'_{ii},$$

$$(b) \quad |D'_m D_m| = 2^{m(m-1)/2}.$$

8.66 Prove the results of Theorem 8.37.

8.67 If A is an $m \times m$ matrix, show that

$$(a) \quad D_m D_m^+ (A \otimes A) D_m = (A \otimes A) D_m,$$

$$(b) \quad \{D_m^+ (A \otimes A) D_m\}^i = D_m^+ (A^i \otimes A^i) D_m, \text{ where } i \text{ is any positive integer.}$$

8.68 Let A be an $m \times m$ nonsingular symmetric matrix and α be a scalar. Show that

$$\begin{aligned} & (D'_m \{A \otimes A + \alpha \text{vec}(A) \text{vec}(A)'\} D_m)^{-1} \\ &= D_m^+ \{A^{-1} \otimes A^{-1} - \beta \text{vec}(A^{-1}) \text{vec}(A^{-1})'\} D_m^+, \end{aligned}$$

where $\beta = \alpha/(1 + m\alpha)$.

8.69 If \mathbf{u}_{ij} and E_{ij} are defined as in Problem 8.63, show that

$$L'_m = \sum_{i \geq j} \{\text{vec}(E_{ij})\} \mathbf{u}'_{ij};$$

that is, verify that

$$\sum_{i \geq j} \{\text{vec}(E_{ij})\} \mathbf{u}'_{ij} \mathbf{v}(A) = \text{vec}(A),$$

where A is an arbitrary $m \times m$ lower triangular matrix.

8.70 Prove the following results:

$$(a) \quad L'_m L_m = \sum_{i \geq j} (E_{jj} \otimes E_{ii}), \text{ where } E_{ii} \text{ was defined in Problem 8.63.}$$

(b) If A and B are $m \times m$ lower triangular matrices, then

$$L'_m L_m (A' \otimes B) L'_m = (A' \otimes B) L'_m.$$

(c) If A is an $m \times m$ nonsingular lower triangular matrix, α is a scalar, and $\beta = \alpha/(1 + m\alpha)$, then

$$\begin{aligned} & (L_m \{A' \otimes A + \alpha \text{vec}(A) \text{vec}(A')'\} L'_m)^{-1} \\ &= L_m \{A^{-1'} \otimes A^{-1} - \beta \text{vec}(A^{-1}) \text{vec}(A^{-1'})'\} L'_m. \end{aligned}$$

8.71 Prove Theorem 8.38.

8.72 For $i = 2, \dots, m, j = 1, \dots, i-1$, define the $m(m-1)/2 \times 1$ vector $\tilde{\mathbf{u}}_{ij}$ to be the vector with one in its $\{(j-1)m + i - j(j+1)/2\}$ th position and zeros elsewhere. It can be easily verified that these vectors are the columns of the identity matrix of order $m(m-1)/2$; that is,

$$I_{m(m-1)/2} = (\tilde{\mathbf{u}}_{21}, \dots, \tilde{\mathbf{u}}_{m1}, \tilde{\mathbf{u}}_{32}, \dots, \tilde{\mathbf{u}}_{m2}, \tilde{\mathbf{u}}_{43}, \dots, \tilde{\mathbf{u}}_{m,m-1}).$$

Show that $\tilde{L}'_m = \sum_{i>j} \{\text{vec}(E_{ij})\} \tilde{\mathbf{u}}'_{ij}$; that is, verify that

$$\sum_{i>j} \{\text{vec}(E_{ij})\} \tilde{\mathbf{u}}'_{ij} \tilde{\mathbf{v}}(A) = \text{vec}(A),$$

where A is an arbitrary $m \times m$ strictly lower triangular matrix.

8.73 Prove the results of Theorem 8.39.

8.74 Find a 2×2 nonnegative matrix A that has its spectral radius equal to 1, yet A^k does not converge to anything as $k \rightarrow \infty$.

8.75 Show that the inverse of a nonsingular positive matrix cannot be nonnegative. Show that the inverse of a nonsingular nonnegative matrix A can be nonnegative only if A has exactly one nonzero element in each column.

8.76 Show that if A is a nonnegative matrix and, for some positive integer k , A^k is a positive matrix, then $\rho(A) > 0$.

8.77 It can be shown (see, for example, Horn and Johnson, 2013) that if A is an $m \times m$ nonnegative matrix, then $\rho(A)$ is an eigenvalue of A and a nonnegative eigenvector \mathbf{x} corresponding to the eigenvalue $\rho(A)$ exists. This result is weaker than the result for irreducible nonnegative matrices. For each of the following, find a 2×2 nonnull reducible matrix A , such that the stated condition holds.

(a) $\rho(A) = 0$.

(b) \mathbf{x} is not positive for any \mathbf{x} satisfying $A\mathbf{x} = \rho(A)\mathbf{x}$.

(c) $\rho(A)$ is a multiple eigenvalue.

- 8.78** Verify that the absolute value of each of the eigenvalues of the 2×2 irreducible matrix

$$A = \begin{bmatrix} 0 & 1 \\ 1 & 0 \end{bmatrix}$$

is equal to $\rho(A)$.

- 8.79** We have seen in Theorem 8.41 that $\rho(A)$ is an eigenvalue of A if A is positive. In this exercise, we will use the extension of this result that says $\rho(A)$ is an eigenvalue of A if A is nonnegative (see Horn and Johnson, 2013). In the following, assume that A is an $m \times m$ nonnegative matrix.

- (a) Show that $\rho(I_m + A) = 1 + \rho(A)$.
- (b) Show that if $A^k > (0)$ for some positive integer k , then $\rho(A)$ is a simple eigenvalue of A .
- (c) Apply part (b) on the matrix $(I_m + A)$ to prove Theorem 8.49; that is, prove that for any irreducible nonnegative matrix A , $\rho(A)$ must be a simple eigenvalue.

- 8.80** Consider the homogeneous Markov chain that has three states and the matrix of transition probabilities given by

$$P = \begin{bmatrix} 0.50 & 0.25 & 0 \\ 0.50 & 0.50 & 0.25 \\ 0 & 0.25 & 0.75 \end{bmatrix}.$$

- (a) Show that P is primitive.
- (b) Determine the equilibrium distribution; that is, find the vector π such that $\lim_{t \rightarrow \infty} P^{(t)} = \pi$.

- 8.81** An $m \times m$ matrix A is said to be completely positive (see, for example, Berman and Shaked-Monderer, 2003) if it can be expressed as $A = BB'$ for some $m \times r$ nonnegative matrix B , where $r \leq m$. Clearly every completely positive matrix must be a nonnegative definite matrix and a nonnegative matrix. Show that in the 2×2 case, these conditions are also sufficient; that is, show that a 2×2 matrix A that is nonnegative definite and nonnegative will also be completely positive.

- 8.82** Let A be the $m \times m$ circulant matrix $\text{circ}(a_1, \dots, a_m)$.

- (a) Find the trace of A .
- (b) Find the determinant of A .

- 8.83** Show that the conjugate transpose of the matrix F given in Theorem 8.56 is

$$F^* = \frac{1}{\sqrt{m}} \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ 1 & \theta^{-1} & \theta^{-2} & \cdots & \theta^{-(m-1)} \\ 1 & \theta^{-2} & \theta^{-4} & \cdots & \theta^{-2(m-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \theta^{-(m-1)} & \theta^{-2(m-1)} & \cdots & \theta^{-(m-1)(m-1)} \end{bmatrix}.$$

Then use the geometric series partial sum formula

$$\sum_{j=0}^n r^j = \frac{1 - r^{n+1}}{1 - r}$$

to prove that $F^{-1} = F^*$.

8.84 Let F be defined as in Theorem 8.56, and let $\Gamma = (e_1, e_m, e_{m-1}, \dots, e_2)$. Show that

(a) $F^2 = \Gamma$,

(b) $F^4 = I_m$,

(c) $F^3 = F^*$.

8.85 Let Π_m be the circulant matrix defined in Section 8.9. Show that

(a) $\Pi_m^{m-1} = \Pi_m^{-1}$,

(b) $\Pi_m^m = I_m$,

(c) $\Pi_m^{mn+r} = \Pi_m^r$, for any integers n and r .

8.86 If $A = \text{circ}(a_1, \dots, a_m)$ and $B = \text{circ}(b_1, \dots, b_m)$, find the eigenvalues of $A + B$ and AB .

8.87 Find the eigenvalues of the circulant matrix $A = \text{circ}(1, \dots, 1)$ by using Theorem 8.57.

8.88 Show that if A is a singular circulant matrix, then its Moore–Penrose inverse, A^+ , is also a circulant matrix.

8.89 Find square matrices A and B of the same order, such that A and B are not circulant matrices, yet their product AB is a circulant matrix.

8.90 Let B be the $m \times m$ Jordan block matrix $J_m(0)$. Show that an $m \times m$ matrix A is a Toeplitz matrix if and only if it can be written in the form

$$A = a_0 I_m + \sum_{j=1}^{m-1} (a_j B^j + a_{-j} B^{j'}).$$

8.91 Consider the $m \times m$ Toeplitz matrix

$$A = \begin{bmatrix} 1 & b & b^2 & \cdots & b^{m-1} \\ a & 1 & b & \cdots & b^{m-2} \\ a^2 & a & 1 & \cdots & b^{m-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a^{m-1} & a^{m-2} & a^{m-3} & \cdots & 1 \end{bmatrix},$$

where $ab \neq 1$. Verify by multiplication that the inverse of A is given by

$$A^{-1} = \begin{bmatrix} c & -bc & 0 & \cdots & 0 & 0 \\ -ac & (ab+1)c & -bc & \cdots & 0 & 0 \\ 0 & -ac & (ab+1)c & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & (ab+1)c & -bc \\ 0 & 0 & 0 & \cdots & -ac & c \end{bmatrix},$$

where $c = (1 - ab)^{-1}$. Show that A is singular if $ab = 1$.

- 8.92** Suppose that z_1, \dots, z_{m+1} are independent random variables each having mean 0 and variance 1. Let \mathbf{x} be the $m \times 1$ random vector that has as its i th component,

$$x_i = z_{i+1} - \rho z_i,$$

where ρ is a constant. Show that the covariance matrix of \mathbf{x} is a Toeplitz matrix of the form given in (8.24), and find the values of a_0 and a_1 .

- 8.93** Consider the $m \times m$ symmetric Toeplitz matrix given by

$$A = \begin{bmatrix} 1 & \rho & \rho^2 & \cdots & \rho^{m-1} \\ \rho & 1 & \rho & \cdots & \rho^{m-2} \\ \rho^2 & \rho & 1 & \cdots & \rho^{m-3} \\ \vdots & \vdots & \vdots & & \vdots \\ \rho^{m-1} & \rho^{m-2} & \rho^{m-3} & \cdots & 1 \end{bmatrix},$$

where $0 < \rho < 1$. Use Theorem 7.6 to show that A is positive definite.

- 8.94** Find a Hadamard matrix of order 8.
- 8.95** Give a Hadamard matrix of order 12, thereby illustrating the existence of a Hadamard matrix of order m , where $m \neq 2^n$ for any positive integer n .
- 8.96** Show that the determinant of a Hadamard matrix attains the upper bound of the Hadamard inequality given in Corollary 8.18.1.
- 8.97** Let A, B, C , and D be $m \times m$ matrices with all of their elements equal to $+1$ and -1 , and define H as

$$H = \begin{bmatrix} A & B & C & D \\ -B & A & -D & C \\ -C & D & A & -B \\ -D & -C & B & A \end{bmatrix}.$$

Show that if

$$AA' + BB' + CC' + DD' = 4mI_m$$

and

$$XY' = YX'$$

for every pair of matrices X and Y chosen from A , B , C , and D , then H is a Hadamard matrix of order $4m$.

- 8.98** Let H and E be Hadamard matrices of order $4m$ and $4n$, respectively. Partition these matrices as

$$H = \begin{bmatrix} P & Q \\ R & S \end{bmatrix}, \quad E = \begin{bmatrix} K & L \\ M & N \end{bmatrix},$$

where each submatrix of H is $2m \times 2m$ and each submatrix of E is $2n \times 2n$. Show that if F is defined as

$$\begin{bmatrix} (P+Q) \otimes K + (P-Q) \otimes M & (P+Q) \otimes L + (P-Q) \otimes N \\ (R+S) \otimes K + (R-S) \otimes M & (R+S) \otimes L + (R-S) \otimes N \end{bmatrix},$$

then $\frac{1}{2}F$ is a Hadamard matrix of order $8mn$.

- 8.99** Show that when $m = n$, the Vandermonde matrix A given in (8.27) is non-singular if and only if the m elements of the second row are distinct.
- 8.100** Let A be the Vandermonde matrix given in (8.27) with $m = n$. Prove that if there are r distinct values in the set $\{a_1, \dots, a_m\}$, then $\text{rank}(A) = r$.
- 8.101** Let P be the $m \times m$ orthogonal matrix $(e_m, e_{m-1}, \dots, e_1)$. Show that if A is an $m \times m$ Vandermonde matrix, then PAA' and $AA'P$ are Toeplitz matrices.

9

MATRIX DERIVATIVES AND RELATED TOPICS

9.1 INTRODUCTION

Differential calculus has widespread applications in statistics. For example, estimation procedures such as the maximum likelihood method and the method of least squares use the optimization properties of derivatives, whereas the so-called delta method for obtaining the asymptotic distribution of a function of random variables uses the first derivative to obtain a first-order Taylor series approximation. These and other applications of differential calculus often involve vectors or matrices. In this chapter, we obtain some of the most commonly encountered matrix derivatives.

9.2 MULTIVARIABLE DIFFERENTIAL CALCULUS

We will begin with a brief review of some of the basic notation, concepts, and results of elementary and multivariable differential calculus. Throughout this section, we will assume differentiability or multiple differentiability of the functions we discuss. For more details on the conditions for differentiability, see Magnus and Neudecker (1999). If f is a real-valued function of one variable, x , then its derivative at x , if it exists, is given by

$$f^{(1)}(x) = f'(x) = \frac{d}{dx} f(x) = \lim_{u \rightarrow 0} \frac{f(x+u) - f(x)}{u}.$$

Equivalently, $f'(x)$ is the quantity that gives the first-order Taylor formula for $f(x + u)$. In other words,

$$f(x + u) = f(x) + uf'(x) + r_1(u, x), \quad (9.1)$$

where the remainder $r_1(u, x)$ is a function of u and x satisfying

$$\lim_{u \rightarrow 0} \frac{r_1(u, x)}{u} = 0.$$

The quantity

$$d_u f(x) = uf'(x) \quad (9.2)$$

appearing in (9.1) is called the first differential of f at x with increment u . This increment u is the differential of x . Later we will use dx in place of u , that is, write $f(x + dx)$ instead of $f(x + u)$, to emphasize the fact that u is the differential of x . For notational convenience, we will often denote the differential given in (9.2) simply by df . Generalizations of (9.1) can be obtained by taking higher ordered derivatives; that is, with the i th derivative of f at x defined as

$$f^{(i)}(x) = \frac{d^i}{dx^i} f(x) = \lim_{u \rightarrow 0} \frac{f^{(i-1)}(x + u) - f^{(i-1)}(x)}{u},$$

we have the k th-order Taylor formula

$$\begin{aligned} f(x + u) &= f(x) + \sum_{i=1}^k \frac{u^i f^{(i)}(x)}{i!} + r_k(u, x) \\ &= f(x) + \sum_{i=1}^k \frac{d_u^i f(x)}{i!} + r_k(u, x), \end{aligned}$$

where $r_k(u, x)$ is a function of u and x satisfying

$$\lim_{u \rightarrow 0} \frac{r_k(u, x)}{u^k} = 0,$$

and

$$d_u^i f(x) = u^i f^{(i)}(x),$$

or simply $d^i f$, is the i th differential of f at x with increment u .

The chain rule is a useful formula for calculating the derivative of a composite function. If y , g , and f are functions such that $y(x) = g(f(x))$, then

$$y'(x) = g'(f(x))f'(x). \quad (9.3)$$

If f is a real-valued function of the $n \times 1$ vector $\mathbf{x} = (x_1, \dots, x_n)'$, then its derivative at \mathbf{x} , if it exists, is given by the $1 \times n$ row vector

$$\frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) = \left[\frac{\partial}{\partial x_1} f(\mathbf{x}) \quad \cdots \quad \frac{\partial}{\partial x_n} f(\mathbf{x}) \right],$$

where

$$\frac{\partial}{\partial x_i} f(\mathbf{x}) = \lim_{u_i \rightarrow 0} \frac{f(\mathbf{x} + u_i \mathbf{e}_i) - f(\mathbf{x})}{u_i}$$

is the partial derivative of f with respect to x_i , and \mathbf{e}_i is the i th column of I_n . The first-order Taylor formula analogous to (9.1) is given by

$$f(\mathbf{x} + \mathbf{u}) = f(\mathbf{x}) + \left(\frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) \right) \mathbf{u} + r_1(\mathbf{u}, \mathbf{x}), \quad (9.4)$$

where the remainder, $r_1(\mathbf{u}, \mathbf{x})$, satisfies

$$\lim_{\mathbf{u} \rightarrow \mathbf{0}} \frac{r_1(\mathbf{u}, \mathbf{x})}{(\mathbf{u}'\mathbf{u})^{1/2}} = 0.$$

The second term on the right-hand side of (9.4) is the first differential of f at \mathbf{x} with incremental vector \mathbf{u} ; that is,

$$df = d_{\mathbf{u}} f(\mathbf{x}) = \left(\frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) \right) \mathbf{u} = \sum_{i=1}^n u_i \frac{\partial}{\partial x_i} f(\mathbf{x}).$$

It is important to note the relationship between the first differential and the first derivative; the first differential of f at \mathbf{x} in \mathbf{u} is the first derivative of f at \mathbf{x} times \mathbf{u} . The higher order differentials of f at \mathbf{x} in the vector \mathbf{u} are given by

$$d^i f = d_{\mathbf{u}}^i f(\mathbf{x}) = \sum_{j_1=1}^n \cdots \sum_{j_i=1}^n u_{j_1} \cdots u_{j_i} \frac{\partial^i}{\partial x_{j_1} \cdots \partial x_{j_i}} f(\mathbf{x}),$$

and these differentials appear in the k th-order Taylor formula,

$$f(\mathbf{x} + \mathbf{u}) = f(\mathbf{x}) + \sum_{i=1}^k \frac{d^i f}{i!} + r_k(\mathbf{u}, \mathbf{x}),$$

where the remainder $r_k(\mathbf{u}, \mathbf{x})$ satisfies

$$\lim_{\mathbf{u} \rightarrow \mathbf{0}} \frac{r_k(\mathbf{u}, \mathbf{x})}{(\mathbf{u}'\mathbf{u})^{k/2}} = 0.$$

The second differential, $d^2 f$, can be written as a quadratic form in the vector \mathbf{u} ; that is,

$$d^2 f = \mathbf{u}' H_f \mathbf{u},$$

where H_f , called the Hessian, is the matrix of second-order partial derivatives given by

$$H_f = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} f(\mathbf{x}) & \frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(\mathbf{x}) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_2^2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_n \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_n^2} f(\mathbf{x}) \end{bmatrix}.$$

9.3 VECTOR AND MATRIX FUNCTIONS

Suppose now that f_1, \dots, f_m each is a function of the same $n \times 1$ vector $\mathbf{x} = (x_1, \dots, x_n)'$. These m functions can be conveniently expressed as components of the vector function

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} f_1(\mathbf{x}) \\ \vdots \\ f_m(\mathbf{x}) \end{bmatrix}.$$

The function \mathbf{f} is differentiable at \mathbf{x} if and only if each component function f_i is differentiable at \mathbf{x} . The Taylor formulas from the previous section can be applied component-wise to \mathbf{f} . For instance, the first-order Taylor formula is given by

$$\begin{aligned} \mathbf{f}(\mathbf{x} + \mathbf{u}) &= \mathbf{f}(\mathbf{x}) + \left(\frac{\partial}{\partial \mathbf{x}'} \mathbf{f}(\mathbf{x}) \right) \mathbf{u} + \mathbf{r}_1(\mathbf{u}, \mathbf{x}) \\ &= \mathbf{f}(\mathbf{x}) + d\mathbf{f}(\mathbf{x}) + \mathbf{r}_1(\mathbf{u}, \mathbf{x}), \end{aligned}$$

where the vector remainder $\mathbf{r}_1(\mathbf{u}, \mathbf{x})$ satisfies

$$\lim_{\mathbf{u} \rightarrow \mathbf{0}} \frac{\mathbf{r}_1(\mathbf{u}, \mathbf{x})}{(\mathbf{u}'\mathbf{u})^{1/2}} = \mathbf{0}$$

and the first derivative of \mathbf{f} at \mathbf{x} is given by the $m \times n$ matrix

$$\frac{\partial}{\partial \mathbf{x}'} \mathbf{f}(\mathbf{x}) = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(\mathbf{x}) & \frac{\partial}{\partial x_2} f_1(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_1(\mathbf{x}) \\ \frac{\partial}{\partial x_1} f_2(\mathbf{x}) & \frac{\partial}{\partial x_2} f_2(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_2(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial}{\partial x_1} f_m(\mathbf{x}) & \frac{\partial}{\partial x_2} f_m(\mathbf{x}) & \cdots & \frac{\partial}{\partial x_n} f_m(\mathbf{x}) \end{bmatrix}.$$

This matrix of partial derivatives is sometimes referred to as the Jacobian matrix of \mathbf{f} at \mathbf{x} . Again, it is crucial to understand the relationship between the first differential and the first derivative. If we obtain the first differential of \mathbf{f} at \mathbf{x} in \mathbf{u} and write it in the form

$$d\mathbf{f} = B\mathbf{u},$$

then the $m \times n$ matrix B must be the derivative of \mathbf{f} at \mathbf{x} .

If y and g are real-valued functions satisfying $y(\mathbf{x}) = g(\mathbf{f}(\mathbf{x}))$, then the generalization of the chain rule given in (9.3) is

$$\begin{aligned}\frac{\partial}{\partial x_i} y(\mathbf{x}) &= \sum_{j=1}^m \left(\frac{\partial}{\partial f_j} g(\mathbf{f}) \right) \left(\frac{\partial}{\partial x_i} f_j(\mathbf{x}) \right) \\ &= \left(\frac{\partial}{\partial \mathbf{f}'} g(\mathbf{f}) \right) \left(\frac{\partial}{\partial x_i} \mathbf{f}(\mathbf{x}) \right)\end{aligned}$$

for $i = 1, \dots, n$, or simply

$$\frac{\partial}{\partial \mathbf{x}'} y(\mathbf{x}) = \left(\frac{\partial}{\partial \mathbf{f}'} g(\mathbf{f}) \right) \left(\frac{\partial}{\partial \mathbf{x}'} \mathbf{f}(\mathbf{x}) \right).$$

In some applications, the f_j 's or the x_i 's are arranged in a matrix instead of a vector. Thus, the most general case involves the $p \times q$ matrix function

$$F(X) = \begin{bmatrix} f_{11}(X) & f_{12}(X) & \cdots & f_{1q}(X) \\ f_{21}(X) & f_{22}(X) & \cdots & f_{2q}(X) \\ \vdots & \vdots & & \vdots \\ f_{p1}(X) & f_{p2}(X) & \cdots & f_{pq}(X) \end{bmatrix}$$

of the $m \times n$ matrix X . Results for the vector function $\mathbf{f}(\mathbf{x})$ can be easily extended to the matrix function $F(X)$ by using the vec operator; that is, let \mathbf{f} be the $pq \times 1$ vector function such that $\mathbf{f}(\text{vec}(X)) = \text{vec}(F(X))$. Then, for instance, the Jacobian matrix of F at X is given by the $pq \times mn$ matrix

$$\frac{\partial}{\partial \text{vec}(X)'} \mathbf{f}(\text{vec}(X)) = \frac{\partial}{\partial \text{vec}(X)'} \text{vec}(F(X)),$$

which has as its (i, j) th element the partial derivative of the i th element of $\text{vec}(F(X))$ with respect to the j th element of $\text{vec}(X)$. This could then be used to obtain the first-order Taylor formula for $\text{vec}(F(X + U))$. The differentials of the matrix $F(X)$ are defined by the equations

$$\text{vec}(d^i F) = \text{vec}(d_U^i F(X)) = d^i \mathbf{f} = d_{\text{vec}(U)}^i \mathbf{f}(\text{vec}(X));$$

that is, $d^i F$, the i th-order differential of F at X in the incremental matrix U , is defined as the $p \times q$ matrix obtained by unstacking the i th-order differential of \mathbf{f} at $\text{vec}(X)$ in the incremental vector $\text{vec}(U)$.

Basic properties of vector and matrix differentials follow in a fairly straightforward fashion from the corresponding properties of scalar differentials. We will summarize some of these properties here. If x and y are functions and α is a constant, then the differential operator, d , satisfies

- (a) $d\alpha = 0$,
- (b) $d(\alpha x) = \alpha dx$,
- (c) $d(x + y) = dx + dy$,
- (d) $d(xy) = (dx)y + x(dy)$,
- (e) $dx^\alpha = \alpha x^{\alpha-1} dx$,
- (f) $de^x = e^x dx$,
- (g) $d\log(x) = x^{-1} dx$.

For instance, to illustrate property (d), note that

$$(x + dx)(y + dy) = xy + x(dy) + (dx)y + (dx)(dy),$$

and $d(xy)$ will be given by the first-degree term in dx and dy , which is $(dx)y + x(dy)$ as required. Using these properties and the definition of a matrix differential, it is easily shown that if X and Y are matrix functions and A is a matrix of constants, then

- (h) $dA = (0)$,
- (i) $d(\alpha X) = \alpha dX$,
- (j) $d(X') = (dX)'$,
- (k) $d(X + Y) = dX + dY$,
- (l) $d(XY) = (dX)Y + X(dY)$,
- (m) $d\operatorname{tr}(X) = \operatorname{tr}(dX)$,
- (n) $d\operatorname{vec}(X) = \operatorname{vec}(dX)$,
- (o) $d(X \otimes Y) = (dX) \otimes Y + X \otimes (dY)$,
- (p) $d(X \odot Y) = (dX) \odot Y + X \odot (dY)$.

We will verify property (l). Thus, we must show that the (i, j) th element of the matrix on the left-hand side of the equation, $(d(XY))_{ij}$, is the same as the (i, j) th element on the right-hand side, $(dX)_{i\cdot}(Y)_{\cdot j} + (X)_{i\cdot}(dY)_{\cdot j}$, where X is $m \times n$ and Y is $n \times m$. Using properties (c) and (d), we find that

$$\begin{aligned} (d(XY))_{ij} &= d\{(X)_{i\cdot}(Y)_{\cdot j}\} = d\left\{\sum_{k=1}^n x_{ik}y_{kj}\right\} \\ &= \sum_{k=1}^n d(x_{ik}y_{kj}) = \sum_{k=1}^n \{(dx_{ik})y_{kj} + x_{ik}dy_{kj}\} \\ &= \sum_{k=1}^n (dx_{ik})y_{kj} + \sum_{k=1}^n x_{ik}dy_{kj} \\ &= (dX)_{i\cdot}(Y)_{\cdot j} + (X)_{i\cdot}(dY)_{\cdot j}, \end{aligned}$$

and so (l) is proven.

We illustrate the application of some of these properties first by finding the derivatives of some simple scalar functions of a vector \mathbf{x} , and then by finding the derivatives of some simple matrix functions of a matrix X .

Example 9.1 Let \mathbf{x} be an $m \times 1$ vector of unrelated variables, and define the functions

$$f(\mathbf{x}) = \mathbf{a}'\mathbf{x},$$

where \mathbf{a} is an $m \times 1$ vector of constants, and

$$g(\mathbf{x}) = \mathbf{x}'A\mathbf{x},$$

where A is an $m \times m$ symmetric matrix of constants. The h th component of the $1 \times m$ row vector $\partial f / \partial \mathbf{x}'$ is $\partial f / \partial x_h$ and

$$\frac{\partial}{\partial x_h} f = \frac{\partial}{\partial x_h} \sum_{i=1}^m a_i x_i = \sum_{i=1}^m a_i \left(\frac{\partial}{\partial x_h} x_i \right) = a_h,$$

because

$$\frac{\partial}{\partial x_h} x_i = \begin{cases} 1, & \text{if } i = h, \\ 0, & \text{if } i \neq h. \end{cases}$$

This then implies that

$$\frac{\partial}{\partial \mathbf{x}'} f = \mathbf{a}'.$$

In a similar fashion, we compute the h th component of the $1 \times m$ row vector $\partial g / \partial \mathbf{x}'$ as

$$\begin{aligned} \frac{\partial}{\partial x_h} g &= \frac{\partial}{\partial x_h} \sum_{i=1}^m \sum_{j=1}^m a_{ij} x_i x_j = \sum_{i=1}^m \sum_{j=1}^m a_{ij} \left(\frac{\partial}{\partial x_h} x_i x_j \right) \\ &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} \left\{ \left(\frac{\partial}{\partial x_h} x_i \right) x_j + x_i \left(\frac{\partial}{\partial x_h} x_j \right) \right\} \\ &= \sum_{j=1}^m a_{hj} x_j + \sum_{i=1}^m a_{ih} x_i \\ &= \sum_{j=1}^m a_{jh} x_j + \sum_{i=1}^m a_{ih} x_i = 2 \sum_{i=1}^m a_{ih} x_i, \end{aligned}$$

because $a_{jh} = a_{hj}$. Note that this derivative can be written as $2\mathbf{x}'(A)_{\cdot h}$, so

$$\frac{\partial}{\partial \mathbf{x}'} g = 2\mathbf{x}'A.$$

An alternative approach to computing these derivatives, one that we will use in most of our examples, involves direct calculation of the differential. As a by-product, we obtain the derivative. For instance, the differential of the first function is

$$df = d(\mathbf{a}'\mathbf{x}) = \mathbf{a}'d\mathbf{x}.$$

Since this differential and the derivative are related through the equation

$$df = \left(\frac{\partial}{\partial \mathbf{x}'} f \right) d\mathbf{x},$$

we immediately observe that the derivative is given by

$$\frac{\partial}{\partial \mathbf{x}'} f = \mathbf{a}'.$$

The differential of our second function is given by

$$\begin{aligned} dg &= d(\mathbf{x}'A\mathbf{x}) = d(\mathbf{x}')A\mathbf{x} + \mathbf{x}'d(A\mathbf{x}) = (d\mathbf{x})'A\mathbf{x} + \mathbf{x}'Ad\mathbf{x} \\ &= \{(\mathbf{x}')'A\mathbf{x}\}' + \mathbf{x}'Ad\mathbf{x} = \mathbf{x}'A'd\mathbf{x} + \mathbf{x}'Ad\mathbf{x} = 2\mathbf{x}'Ad\mathbf{x}. \end{aligned}$$

By again making use of the relationship between the differential and the derivative, we observe that

$$\frac{\partial}{\partial \mathbf{x}'} g = 2\mathbf{x}'A.$$

Example 9.2 Let X be an $m \times n$ matrix of unrelated variables, and define the functions

$$F(X) = AX,$$

where A is a $p \times m$ matrix of constants, and

$$G(X) = (X - C)'B(X - C),$$

where B is an $m \times m$ symmetric matrix of constants and C is an $m \times n$ matrix of constants. We will find the Jacobian matrices by first obtaining the differentials of these functions. The derivatives can be obtained from these differentials because, for instance, if we obtain

$$d \operatorname{vec}(F) = W d \operatorname{vec}(X),$$

then the matrix W will be the derivative of $\operatorname{vec}(F(X))$ with respect to $\operatorname{vec}(X)$. For our first function, we find that

$$dF = d(AX) = AdX,$$

so that

$$\begin{aligned} d \operatorname{vec}(F) &= \operatorname{vec}(dF) = \operatorname{vec}(AdX) = (I_n \otimes A) \operatorname{vec}(dX) \\ &= (I_n \otimes A) d \operatorname{vec}(X), \end{aligned}$$

where we have used Theorem 8.11. Thus, we must have

$$\frac{\partial}{\partial \operatorname{vec}(X)'} \operatorname{vec}(F) = I_n \otimes A.$$

The differential of our second function is

$$\begin{aligned} dG &= d\{(X - C)'B(X - C)\} \\ &= \{d(X' - C')\}B(X - C) + (X - C)'B\{d(X - C)\} \\ &= (dX)'B(X - C) + (X - C)'BdX, \end{aligned}$$

from which we obtain

$$\begin{aligned} d \operatorname{vec}(G) &= \{(X - C)'B \otimes I_n\} \operatorname{vec}(dX') + \{I_n \otimes (X - C)'B\} \operatorname{vec}(dX) \\ &= \{(X - C)'B \otimes I_n\} K_{nn} \operatorname{vec}(dX) + \{I_n \otimes (X - C)'B\} \operatorname{vec}(dX) \\ &= K_{nn} \{I_n \otimes (X - C)'B\} \operatorname{vec}(dX) + \{I_n \otimes (X - C)'B\} \operatorname{vec}(dX) \\ &= (I_{n^2} + K_{nn}) \{I_n \otimes (X - C)'B\} \operatorname{vec}(dX) \\ &= 2N_n \{I_n \otimes (X - C)'B\} d \operatorname{vec}(X), \end{aligned}$$

where we have used properties of the vec operator and the commutation matrix developed in Chapter 8. Consequently, we have

$$\frac{\partial}{\partial \operatorname{vec}(X)'} \operatorname{vec}(G) = 2N_n \{I_n \otimes (X - C)'B\}.$$

In Example 9.3, we show how we can use the Jacobian matrix of the simple transformation $\mathbf{z} = \mathbf{c} + A\mathbf{x}$ to obtain the multivariate normal density function given in (1.13).

Example 9.3 Suppose that \mathbf{z} is an $m \times 1$ random vector with density function $f_1(\mathbf{z})$ that is positive for all $\mathbf{z} \in S_1 \subseteq R^m$. Let the $m \times 1$ vector $\mathbf{x} = \mathbf{x}(\mathbf{z})$ represent a one-to-one transformation of S_1 onto $S_2 \subseteq R^m$, so that the inverse transformation $\mathbf{z} = \mathbf{z}(\mathbf{x})$, $\mathbf{x} \in S_2$ is unique. Denote the Jacobian matrix of \mathbf{z} at \mathbf{x} as

$$J = \frac{\partial}{\partial \mathbf{x}'} \mathbf{z}(\mathbf{x}).$$

If the partial derivatives in J exist and are continuous functions on the set S_2 , then the density of \mathbf{x} is given by

$$f_2(\mathbf{x}) = f_1(\mathbf{z}(\mathbf{x})) |J|.$$

We will use the formula above to obtain the multivariate normal density, given in (1.13), from the standard normal density. Now recall that, by definition, $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ if \mathbf{x} can be expressed as $\mathbf{x} = \boldsymbol{\mu} + T\mathbf{z}$, where $TT' = \Omega$ and the components of \mathbf{z} , z_1, \dots, z_m are independently distributed each as $N(0, 1)$. Thus, the density function of \mathbf{z} is given by

$$f_1(\mathbf{z}) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{1}{2}z_i^2\right) = \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\mathbf{z}'\mathbf{z}\right).$$

The differential of the inverse transformation $\mathbf{z} = T^{-1}(\mathbf{x} - \boldsymbol{\mu})$ is $d\mathbf{z} = T^{-1}d\mathbf{x}$, and so the necessary Jacobian matrix is $J = T^{-1}$. Consequently, we find that the density of \mathbf{x} is given by

$$\begin{aligned} f_2(\mathbf{x}) &= \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2}\{T^{-1}(\mathbf{x} - \boldsymbol{\mu})\}'T^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) |T^{-1}| \\ &= \frac{1}{(2\pi)^{m/2}|T|} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})'T^{-1'}T^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= \frac{1}{(2\pi)^{m/2}|\Omega|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})'\Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\right), \end{aligned}$$

because

$$\Omega^{-1} = (TT')^{-1} = T^{-1'}T^{-1}$$

and

$$|\Omega|^{1/2} = |TT'|^{1/2} = |T|^{1/2}|T'|^{1/2} = |T|^{1/2}|T|^{1/2} = |T|.$$

9.4 SOME USEFUL MATRIX DERIVATIVES

In this section, we will obtain the differentials and the corresponding derivatives of some important scalar functions and matrix functions of matrices. Throughout this section, when dealing with functions of the form $f(X)$ or $F(X)$, we will assume that the $m \times n$ matrix X is composed of mn unrelated variables; that is, X is assumed not to have any particular structure such as symmetry, triangularity, and so on. We begin with some scalar functions of X .

Theorem 9.1 Let X be an $m \times m$ matrix, and let $X_{\#}$ denote the adjoint matrix of X . Then

$$(a) \quad d\{\text{tr}(X)\} = \text{vec}(I_m)'d\text{vec}(X); \quad \frac{\partial}{\partial \text{vec}(X)'}\text{tr}(X) = \text{vec}(I_m)',$$

- (b) $d|X| = \text{tr}(X_{\#}dX)$; $\frac{\partial}{\partial \text{vec}(X)'}|X| = \text{vec}(X'_{\#})'$,
 (c) and if X is nonsingular,
 $d|X| = |X|\text{tr}(X^{-1}dX)$; $\frac{\partial}{\partial \text{vec}(X)'}|X| = |X|\text{vec}(X^{-1})'$.

Proof. Part (a) follows directly from the fact that

$$\begin{aligned} d \text{tr}(X) &= \text{tr}(dX) = \text{tr}(I_m dX) = \text{vec}(I_m)' \text{vec}(dX) \\ &= \text{vec}(I_m)' d \text{vec}(X), \end{aligned}$$

with the third equality following from Theorem 8.10. Since $X_{\#}$ is the transpose of the matrix of cofactors of X , to obtain the derivative in (b), we simply need to show that

$$\frac{\partial}{\partial x_{ij}}|X| = X_{ij},$$

where X_{ij} is the cofactor of x_{ij} . By using the cofactor expansion formula on the i th row of X , we can write the determinant of X as

$$|X| = \sum_{k=1}^m x_{ik} X_{ik}.$$

Note that for each k , X_{ik} is a determinant computed after deleting the i th row so that each X_{ik} does not involve the element x_{ij} . Consequently, we have

$$\frac{\partial}{\partial x_{ij}}|X| = \frac{\partial}{\partial x_{ij}} \sum_{k=1}^m x_{ik} X_{ik} = \sum_{k=1}^m \left(\frac{\partial}{\partial x_{ij}} x_{ik} \right) X_{ik} = X_{ij}.$$

Using the relationship between the first differential and derivative, we then get

$$d|X| = \{\text{vec}(X'_{\#})\}' \text{vec}(dX) = \text{tr}(X_{\#}dX)$$

as required. Part (c) follows directly from (b) because if X is nonsingular, then $X^{-1} = |X|^{-1}X_{\#}$. \square

An immediate consequence of Theorem 9.1(c) is Corollary 9.1.1.

Corollary 9.1.1 Let X be an $m \times m$ nonsingular matrix. Then

$$d\{\log(|X|)\} = \text{tr}(X^{-1}dX); \quad \frac{\partial}{\partial \text{vec}(X)'} \log(|X|) = \text{vec}(X^{-1})'.$$

Theorem 9.2 gives the differential and derivative of the inverse of a nonsingular matrix.

Theorem 9.2 If X is a nonsingular $m \times m$ matrix, then

$$dX^{-1} = -X^{-1}(dX)X^{-1}; \quad \frac{\partial}{\partial \text{vec}(X)'} \text{vec}(X^{-1}) = -(X^{-1'} \otimes X^{-1}).$$

Proof. Computing the differential of both sides of the equation $I_m = XX^{-1}$, we find that

$$(0) = dI_m = d(XX^{-1}) = (dX)X^{-1} + X(dX^{-1}).$$

Premultiplying this equation by X^{-1} and then solving for dX^{-1} yields

$$dX^{-1} = -X^{-1}(dX)X^{-1},$$

which leads to

$$\begin{aligned} d \text{vec}(X^{-1}) &= \text{vec}(dX^{-1}) = -\text{vec}(X^{-1}(dX)X^{-1}) \\ &= -(X^{-1'} \otimes X^{-1})\text{vec}(dX) \\ &= -(X^{-1'} \otimes X^{-1})d \text{vec}(X), \end{aligned}$$

and so the proof is complete. \square

A natural generalization of Theorem 9.2 is one that gives the differential and derivative of the Moore–Penrose inverse of a matrix. Theorem 9.3 gives the form of these when they exist at a matrix X . Conditions for the differentiability of X^+ can be found in Magnus and Neudecker (1999).

Theorem 9.3 If X is an $m \times n$ matrix and X^+ is its Moore–Penrose inverse, then

$$dX^+ = (I_n - X^+X)(dX')X^{+'}X^+ + X^+X^{+'}(dX')(I_m - XX^+) - X^+(dX)X^+$$

and

$$\begin{aligned} \frac{\partial}{\partial \text{vec}(X)'} \text{vec}(X^+) &= \{X^{+'}X^+ \otimes (I_n - X^+X) + (I_m - XX^+) \otimes X^+X^{+'}\} \\ &\quad \times K_{mn} - (X^{+'} \otimes X^+). \end{aligned}$$

Proof. Note that

$$d(XX^+) = (dX)X^+ + XdX^+,$$

from which we get

$$XdX^+ = d(XX^+) - (dX)X^+. \quad (9.5)$$

Since $X^+ = X^+XX^+$, we also have

$$\begin{aligned} dX^+ &= d(X^+XX^+) = d(X^+X)X^+ + X^+XdX^+ \\ &= d(X^+X)X^+ + X^+d(XX^+) - X^+(dX)X^+, \end{aligned} \quad (9.6)$$

where we have used (9.5) in the last step. Thus, if we obtain expressions for $d(X^+X)$ and $d(XX^+)$ in terms of dX , we can then find dX^+ . To find $d(XX^+)$, we use the fact that XX^+ is symmetric and idempotent to get

$$\begin{aligned} d(XX^+) &= d(XX^+XX^+) = d(XX^+)XX^+ + XX^+d(XX^+) \\ &= d(XX^+)XX^+ + (d(XX^+)XX^+)', \end{aligned} \quad (9.7)$$

because $d(XX^+)' = d((XX^+)') = d(XX^+)$. However,

$$d(XX^+)X = dX - XX^+dX = (I_m - XX^+)dX, \quad (9.8)$$

because $X = XX^+X$ implies that

$$dX = d(XX^+X) = d(XX^+)X + XX^+dX.$$

Now substituting (9.8) in (9.7), we find that

$$\begin{aligned} d(XX^+) &= (I_m - XX^+)(dX)X^+ + \{(I_m - XX^+)(dX)X^+\}' \\ &= (I_m - XX^+)(dX)X^+ + X^{+'}(dX')(I_m - XX^+). \end{aligned} \quad (9.9)$$

By using the fact that X^+X is symmetric and idempotent, we can show in a similar fashion that

$$d(X^+X) = X^+(dX)(I_n - X^+X) + (I_n - X^+X)(dX')X^{+'}. \quad (9.10)$$

Substituting (9.9) and (9.10) into (9.6), and noting that $(I_n - X^+X)X^+ = (0)$ and $X^+(I_m - XX^+) = (0)$, we get

$$dX^+ = (I_n - X^+X)(dX')X^{+'}X^+ + X^+X^{+'}(dX')(I_m - XX^+) - X^+(dX)X^+,$$

as is required. When we take the vec of both sides of the equation above, we get

$$\begin{aligned}
 d \text{vec}(X^+) &= \{X^{+'}X^+ \otimes (I_n - X^+X)\} \text{vec}(dX') \\
 &\quad + \{(I_m - XX^+) \otimes X^+X^{+'}\} \text{vec}(dX') \\
 &\quad - (X^{+'} \otimes X^+) \text{vec}(dX) \\
 &= \{X^{+'}X^+ \otimes (I_n - X^+X) \\
 &\quad + (I_m - XX^+) \otimes X^+X^{+'}\} K_{mn} d \text{vec}(X) \\
 &\quad - (X^{+'} \otimes X^+) d \text{vec}(X),
 \end{aligned}$$

which yields the required expression for the differential. \square

9.5 DERIVATIVES OF FUNCTIONS OF PATTERNED MATRICES

In this section, we consider the computation of the derivative of a function of an $m \times n$ matrix X when some of the variables of X are related to one another. In particular, we will focus on the situation in which X is square and symmetric. For a more general treatment of the topic of derivatives of functions of patterned matrices, see Nel (1980).

If X is an $m \times m$ symmetric matrix of variables, then because of the symmetry it only contains $m(m+1)/2$ mathematically independent variables. These variables are precisely the variables comprising the vector $\mathbf{v}(X)$. If $\mathbf{f}(X)$ is some vector function of the matrix X , then the derivative of \mathbf{f} will be given by the matrix

$$\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{f}(X).$$

We can compute derivatives of this form by using the derivative

$$\frac{\partial}{\partial \text{vec}(X)'} \mathbf{f}(X)$$

for a general nonsymmetric matrix X , along with the chain rule. Specifically, from the chain rule we have

$$\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{f}(X) = \left(\frac{\partial}{\partial \text{vec}(X)'} \mathbf{f}(X) \right) \left(\frac{\partial}{\partial \mathbf{v}(X)'} \text{vec}(X) \right).$$

It must be emphasized here that the first of the two derivatives on the right-hand side of this equation is computed ignoring the symmetry of X . The second of these two derivatives can be conveniently expressed if we make use of the duplication matrix D_m . Since $D_m \mathbf{v}(X) = \text{vec}(X)$, we immediately get $D_m d \mathbf{v}(X) = d \text{vec}(X)$, and so

$$\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{f}(X) = \left(\frac{\partial}{\partial \text{vec}(X)'} \mathbf{f}(X) \right) D_m.$$

Consequently, the following results follow directly from Theorem 9.1, Theorem 9.2, and Theorem 9.3.

Theorem 9.4 Let X be an $m \times m$ symmetric matrix of variables. Then

- (a) $\frac{\partial}{\partial \mathbf{v}(X)'} |X| = \text{vec}(X'_{\#})' D_m,$
- (b) $\frac{\partial}{\partial \mathbf{v}(X)'} \text{vec}(X^{-1}) = -(X^{-1} \otimes X^{-1}) D_m$ if X is nonsingular,
- (c) $\frac{\partial}{\partial \mathbf{v}(X)'} \text{vec}(X^+) = (\{X^+ X^+ \otimes (I_m - X^+ X) + (I_m - X X^+) \otimes X^+ X^+\} K_{mm} - (X^+ \otimes X^+)) D_m.$

The derivatives given in (b) and (c) of Theorem 9.4 still have some redundant elements because of the symmetry of X^{-1} and X^+ . In general, if X is an $m \times m$ symmetric matrix of variables and the $m \times m$ matrix function $F(X)$ is also symmetric, then all derivatives of elements of $F(X)$ with respect to elements of X will be contained in the matrix derivative

$$\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{v}\{F(X)\}.$$

This matrix derivative can be easily computed from the derivative

$$A = \frac{\partial}{\partial \mathbf{v}(X)'} \text{vec}\{F(X)\} \quad (9.11)$$

by again using the relationship $\text{vec}(F) = D_m \mathbf{v}(F)$. Thus, because (9.11) implies that $d \text{vec}(F) = A d \mathbf{v}(X)$, we have

$$D_m d \mathbf{v}(F) = A d \mathbf{v}(X)$$

or

$$D_m^+ D_m d \mathbf{v}(F) = d \mathbf{v}(F) = D_m^+ A d \mathbf{v}(X),$$

because $D_m^+ D_m = I_{m(m+1)/2}$ by Theorem 8.32. Using this result, we obtain the derivatives in Corollary 9.4.1.

Corollary 9.4.1 Let X be an $m \times m$ symmetric matrix of variables. Then

- (a) $\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{v}(X^{-1}) = -D_m^+ (X^{-1} \otimes X^{-1}) D_m$ if X is nonsingular,
- (b) $\frac{\partial}{\partial \mathbf{v}(X)'} \mathbf{v}(X^+) = D_m^+ (\{X^+ X^+ \otimes (I_m - X^+ X) + (I_m - X X^+) \otimes X^+ X^+\} K_{mm} - (X^+ \otimes X^+)) D_m.$

9.6 THE PERTURBATION METHOD

The perturbation method is a technique, closely related to the method using the differential operator, for finding successive terms in a Taylor expansion formula. In this section, we will use this method to obtain Taylor formulas for some important matrix functions. A more rigorous treatment of this subject can be found in texts such as Hinch (1991), Kato (1982), or Nayfeh (1981).

Suppose that the elements of dX are small, which we can emphasize by writing $dX = \epsilon Y$, where ϵ is a small scalar and Y is an $m \times n$ matrix. Then $X + \epsilon Y$ represents a small perturbation of the $m \times n$ matrix X . The Taylor formula for the vector function \mathbf{f} of X would then be of the form

$$\mathbf{f}(X + \epsilon Y) = \mathbf{f}(X) + \sum_{i=1}^{\infty} \epsilon^i \mathbf{g}_i(X, Y),$$

where $\mathbf{g}_i(X, Y)$ represents some vector function of the two matrices, X and Y . Similarly, if we have a matrix function F , then the expansion would be of the form

$$F(X + \epsilon Y) = F(X) + \sum_{i=1}^{\infty} \epsilon^i G_i(X, Y). \quad (9.12)$$

Our goal is to determine the first few terms in the summations given above. These then can be applied in an approximation of $\mathbf{f}(X + \epsilon Y)$ or $F(X + \epsilon Y)$ when ϵ is small. For instance, suppose that $m = n$ and our function is the matrix inverse function; that is, $F(X) = X^{-1}$. For notational simplicity, write $G_i(X, Y) = G_i$, and suppose that the $m \times m$ matrices X and $(X + \epsilon Y)$ are nonsingular. Then (9.12) can be written

$$(X + \epsilon Y)^{-1} = X^{-1} + \epsilon G_1 + \epsilon^2 G_2 + \epsilon^3 G_3 + \cdots.$$

However, we must have

$$\begin{aligned} I_m &= (X + \epsilon Y)(X + \epsilon Y)^{-1} \\ &= (X + \epsilon Y)(X^{-1} + \epsilon G_1 + \epsilon^2 G_2 + \epsilon^3 G_3 + \cdots) \\ &= I_m + \epsilon(YX^{-1} + XG_1) + \epsilon^2(YG_1 + XG_2) \\ &\quad + \epsilon^3(YG_2 + XG_3) + \cdots. \end{aligned}$$

If this result is to hold for all ϵ , then we must have $(YX^{-1} + XG_1) = (0)$ or, equivalently,

$$G_1 = -X^{-1}YX^{-1}.$$

Similarly, we must have $(YG_1 + XG_2) = (0)$ so that

$$G_2 = -X^{-1}YG_1 = X^{-1}YX^{-1}YX^{-1},$$

and in fact, it should be apparent that we have the recursive relationship

$$G_h = -X^{-1}YG_{h-1}.$$

As a result, we have

$$\begin{aligned} (X + \epsilon Y)^{-1} &= X^{-1} - \epsilon X^{-1}YX^{-1} + \epsilon^2 X^{-1}YX^{-1}YX^{-1} \\ &\quad - \epsilon^3 X^{-1}YX^{-1}YX^{-1}YX^{-1} + \cdots, \end{aligned}$$

or, if we return to the notation $dX = \epsilon Y$,

$$\begin{aligned} (X + dX)^{-1} &= X^{-1} - X^{-1}(dX)X^{-1} + X^{-1}(dX)X^{-1}(dX)X^{-1} \\ &\quad - X^{-1}(dX)X^{-1}(dX)X^{-1}(dX)X^{-1} + \cdots. \end{aligned}$$

Next, we will use this perturbation method to determine the first few terms in the Taylor series expansion for an eigenvalue of a symmetric matrix. Such an expansion will be possible only if the corresponding eigenvalue of the unperturbed matrix X is distinct. We will first consider the special case in which X is a diagonal matrix.

Theorem 9.5 Suppose $X = \text{diag}(x_1, \dots, x_m)$, where $x_1 \geq \cdots \geq x_{l-1} > x_l > x_{l+1} \geq \cdots \geq x_m$, so that the l th diagonal element x_l differs from the other diagonal elements of X . Let U be an $m \times m$ symmetric matrix, and denote the l th largest eigenvalue and corresponding normalized eigenvector of $X + U$ by $\lambda_l(X + U)$ and $\gamma_l(X + U)$, respectively. Then

$$\begin{aligned} \lambda_l(X + U) &\approx x_l + u_{ll} - \sum_{i \neq l} \frac{u_{il}^2}{(x_i - x_l)} - \sum_{i \neq l} \frac{u_{il}u_{il}^2}{(x_i - x_l)^2} \\ &\quad + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il}u_{jl}u_{ij}}{(x_i - x_l)(x_j - x_l)}, \\ \gamma_{ll}(X + U) &\approx 1 - \frac{1}{2} \sum_{i \neq l} \frac{u_{il}^2}{(x_i - x_l)^2} - \sum_{i \neq l} \frac{u_{il}u_{il}^2}{(x_i - x_l)^3} \\ &\quad + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il}u_{jl}u_{ij}}{(x_i - x_l)^2(x_j - x_l)}, \end{aligned}$$

and for $h \neq l$,

$$\begin{aligned} \gamma_{hl}(X + U) &\approx -\frac{u_{hl}}{(x_h - x_l)} - \frac{u_{ll}u_{hl}}{(x_h - x_l)^2} \\ &\quad + \sum_{i \neq l} \frac{u_{il}u_{hi}}{(x_h - x_l)(x_i - x_l)} - \frac{u_{il}^2u_{hl}}{(x_h - x_l)^3} \end{aligned}$$

$$\begin{aligned}
& + \sum_{i \neq l} \frac{u_{il} u_{hi} u_{il}}{(x_h - x_l)^2 (x_i - x_l)} + \sum_{i \neq l} \frac{u_{hl} u_{il}^2}{(x_h - x_l)^2 (x_i - x_l)} \\
& + \sum_{i \neq l} \frac{u_{il} u_{hi} u_{il}}{(x_h - x_l)(x_i - x_l)^2} \\
& - \sum_{i \neq l} \sum_{j \neq l} \frac{u_{hi} u_{ij} u_{jl}}{(x_h - x_l)(x_i - x_l)(x_j - x_l)} \\
& + \frac{1}{2} \sum_{i \neq l} \frac{u_{hl} u_{il}^2}{(x_h - x_l)(x_i - x_l)^2},
\end{aligned}$$

where $\gamma_{hl}(X + U)$ denotes the h th element of $\gamma_l(X + U)$, and these approximations are accurate up through third-order terms in the u 's.

Proof. Here U is the perturbation matrix, and we wish to write $\lambda_l = \lambda_l(X + U)$ and $\gamma_l = \gamma_l(X + U)$ in the form

$$\lambda_l = x_l + a_1 + a_2 + a_3 + \cdots, \quad (9.13)$$

$$\gamma_l = e_l + b_1 + b_2 + b_3 + \cdots, \quad (9.14)$$

where a_i and b_i only involve i th degree terms in the elements of U . Substituting these expressions in the defining equation $(X + U)\gamma_l = \lambda_l \gamma_l$, and then equating i th degree terms in the elements of U on the left-hand side of this equation to those on the right-hand side, we obtain

$$X e_l = x_l e_l, \quad (9.15)$$

$$X b_1 + U e_l = x_l b_1 + a_1 e_l, \quad (9.16)$$

$$X b_2 + U b_1 = x_l b_2 + a_1 b_1 + a_2 e_l, \quad (9.17)$$

$$X b_3 + U b_2 = x_l b_3 + a_1 b_2 + a_2 b_1 + a_3 e_l. \quad (9.18)$$

In a similar fashion, the normalizing equation $\gamma'_l \gamma_l = 1$ yields the identities

$$e'_l e_l = 1, \quad (9.19)$$

$$e'_l b_1 + b'_1 e_l = 0, \quad (9.20)$$

$$e'_l b_2 + b'_1 b_1 + b'_2 e_l = 0, \quad (9.21)$$

$$e'_l b_3 + b'_1 b_2 + b'_2 b_1 + b'_3 e_l = 0. \quad (9.22)$$

Equations (9.15) and (9.19) are trivially true, whereas equations (9.16) and (9.20) can be used to find a_1 and b_1 . Premultiplying (9.16) by e'_l and then solving for a_1 , we find that

$$a_1 = e'_l U e_l = u_{ll}, \quad (9.23)$$

because $e'_l X \mathbf{b}_1 = x_l e'_l \mathbf{b}_1$ follows from (9.15). We can then rewrite (9.16) as the system of linear equations

$$(X - x_l I_m) \mathbf{b}_1 = -(U - u_{ll} I_m) \mathbf{e}_l,$$

with the general solution for \mathbf{b}_1 given by

$$\mathbf{b}_1 = -(X - x_l I_m)^+ (U - u_{ll} I_m) \mathbf{e}_l + c_1 \mathbf{e}_l,$$

where c_1 is an arbitrary constant. Since $(X - x_l I_m)^+ \mathbf{e}_l = \mathbf{0}$ and (9.20) implies that $e'_l \mathbf{b}_1 = 0$, it follows that $c_1 = 0$, and thus,

$$\mathbf{b}_1 = -(X - x_l I_m)^+ U \mathbf{e}_l. \quad (9.24)$$

Next, we will use (9.17) and (9.21) to find a_2 and \mathbf{b}_2 . Premultiplying (9.17) by e'_l and then solving for a_2 , we find, after we again use the fact that $e'_l \mathbf{b}_1 = 0$, that

$$a_2 = e'_l U \mathbf{b}_1 = -e'_l U (X - x_l I_m)^+ U \mathbf{e}_l. \quad (9.25)$$

Rewriting (9.17) as the system of equations in \mathbf{b}_2 ,

$$(X - x_l I_m) \mathbf{b}_2 = a_2 \mathbf{e}_l - (U - a_1 I_m) \mathbf{b}_1,$$

which for any scalar c_2 has as a solution,

$$\mathbf{b}_2 = (X - x_l I_m)^+ \{a_2 \mathbf{e}_l - (U - a_1 I_m) \mathbf{b}_1\} + c_2 \mathbf{e}_l.$$

Now because $(X - x_l I_m)^+ \mathbf{e}_l = \mathbf{0}$ and (9.21) implies that $e'_l \mathbf{b}_2 = -\frac{1}{2} \mathbf{b}'_1 \mathbf{b}_1$, we find that

$$c_2 = -\frac{1}{2} \mathbf{b}'_1 \mathbf{b}_1 = -\frac{1}{2} e'_l U \{(X - x_l I_m)^+\}^2 U \mathbf{e}_l = -\frac{1}{2} \sum_{i \neq l} \frac{u_{il}^2}{(x_i - x_l)^2},$$

and so with this value for c_2 , the solution for \mathbf{b}_2 is given by

$$\mathbf{b}_2 = (X - x_l I_m)^+ (U - u_{ll} I_m) (X - x_l I_m)^+ U \mathbf{e}_l + c_2 \mathbf{e}_l. \quad (9.26)$$

To find a_3 , premultiply (9.18) by e'_l and solve for a_3 , after using $e'_l \mathbf{b}_1 = 0$, to get

$$\begin{aligned} a_3 &= e'_l (U - a_1 I_m) \mathbf{b}_2 \\ &= e'_l (U - u_{ll} I_m) \{(X - x_l I_m)^+ (U - u_{ll} I_m) (X - x_l I_m)^+ U \mathbf{e}_l + c_2 \mathbf{e}_l\} \\ &= e'_l U (X - x_l I_m)^+ (U - u_{ll} I_m) (X - x_l I_m)^+ U \mathbf{e}_l. \end{aligned} \quad (9.27)$$

Equation (9.18) can be expressed as

$$(X - x_l I_m) \mathbf{b}_3 = a_3 \mathbf{e}_l + a_2 \mathbf{b}_1 - (U - a_1 I_m) \mathbf{b}_2,$$

so that the solution for \mathbf{b}_3 will be given by

$$\begin{aligned} \mathbf{b}_3 &= (X - x_l I_m)^+ \{a_3 \mathbf{e}_l + a_2 \mathbf{b}_1 - (U - a_1 I_m) \mathbf{b}_2\} + c_3 \mathbf{e}_l \\ &= (X - x_l I_m)^+ \{a_2 \mathbf{b}_1 - (U - a_1 I_m) \mathbf{b}_2\} + c_3 \mathbf{e}_l, \end{aligned} \quad (9.28)$$

where c_3 is an arbitrary constant. By premultiplying this equation by \mathbf{e}_l' and using $\mathbf{e}_l' \mathbf{b}_3 = -\mathbf{b}_1' \mathbf{b}_2$ that follows from (9.22), we find that

$$\begin{aligned} c_3 &= -\mathbf{b}_1' \mathbf{b}_2 = \mathbf{e}_l' U \{(X - x_l I_m)^+\}^2 (U - u_{ll} I_m) (X - x_l I_m)^+ U \mathbf{e}_l \\ &= - \sum_{i \neq l} \frac{u_{il} u_{il}^2}{(x_i - x_l)^3} + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il} u_{jl} u_{ij}}{(x_i - x_l)^2 (x_j - x_l)}. \end{aligned}$$

The results now follow by substituting (9.23), (9.25) and (9.27) in (9.13), and (9.24), (9.26) and (9.28) in (9.14). \square

We can use Theorem 9.5 to obtain expansion formulas for a general symmetric matrix; that is, if Z is an $m \times m$ symmetric matrix and W is its associated symmetric perturbation matrix, then we can obtain expansion formulas for $\lambda_l(Z + W)$ and $\gamma_l(Z + W)$. Let $Z = QXQ'$ be the spectral decomposition of Z , so that $X = \text{diag}(x_1, \dots, x_m)$, with x_l being an eigenvalue of Z corresponding to the eigenvector \mathbf{q}_l , which is the l th column of Q . As in Theorem 9.5, we assume that x_l is a distinct eigenvalue. If we let $U = Q'WQ$, then the eigenvalue-eigenvector equation

$$(Z + W)\gamma_l(Z + W) = \lambda_l(Z + W)\gamma_l(Z + W)$$

can be equivalently expressed as

$$(X + U)Q'\gamma_l(Z + W) = \lambda_l(Z + W)Q'\gamma_l(Z + W);$$

that is, U is the perturbation matrix of X and $\lambda_l(Z + W)$ is an eigenvalue of $(X + U)$ corresponding to the eigenvector $Q'\gamma_l(Z + W)$. Thus, if we use the elements of $U = QWQ'$ in place of those of U in the formulas in Theorem 9.5, we will obtain expansions for $\lambda_l(Z + W)$ and $Q'\gamma_l(Z + W)$. For instance, first-order approximations of $\lambda_l(Z + W)$ and $\gamma_l(Z + W)$ are given by

$$\begin{aligned} \lambda_l(Z + W) &\approx x_l + \mathbf{q}_l' W \mathbf{q}_l, \\ \gamma_l(Z + W) &\approx Q \{ \mathbf{e}_l - (X - x_l I_m)^+ (Q' W Q) \mathbf{e}_l \} \\ &= \mathbf{q}_l - (Z - x_l I_m)^+ W \mathbf{q}_l. \end{aligned}$$

Theorem 9.6 is an immediate consequence of these first-order Taylor expansion formulas.

Theorem 9.6 Let $\lambda_l(Z)$ be the eigenvalue function defined on $m \times m$ symmetric matrices Z , and let $\gamma_l(Z)$ be a corresponding normalized eigenvector. If the matrix Z is such that the eigenvalue $\lambda_l(Z)$ is distinct, then differentials and derivatives at that matrix Z are given by

$$\begin{aligned} d\lambda_l &= \gamma'_l(dZ)\gamma_l, & \frac{\partial}{\partial v(Z)'} \lambda_l(Z) &= (\gamma'_l \otimes \gamma'_l)D_m, \\ d\gamma_l &= -(Z - \lambda_l I_m)^+(dZ)\gamma_l, \\ \frac{\partial}{\partial v(Z)'} \gamma_l(Z) &= -\{\gamma'_l \otimes (Z - \lambda_l I_m)^+\}D_m. \end{aligned}$$

The expansions given in and immediately after Theorem 9.5 do not hold when the eigenvalue x_l is not distinct. Suppose, for instance, that again $x_1 \geq \cdots \geq x_m$, but now $x_l = x_{l+1} = \cdots = x_{l+r-1}$, so that the value x_l is repeated as an eigenvalue of $Z = QXQ'$, r times. In this case, we can get expansions for $\bar{\lambda}_{l,l+r-1}(Z + W)$, the average of the perturbed eigenvalues $\lambda_l(Z + W), \dots, \lambda_{l+r-1}(Z + W)$, and the total eigenprojection Φ_l associated with this collection of eigenvalues; if $P_{Z+W}\{\lambda_{l+i-1}(Z + W)\}$ represents the eigenprojection of $Z + W$ associated with the eigenvalue $\lambda_{l+i-1}(Z + W)$, then this total eigenprojection is given by

$$\begin{aligned} \Phi_l &= \sum_{i=1}^r P_{Z+W}\{\lambda_{l+i-1}(Z + W)\} \\ &= \sum_{i=1}^r \gamma_{l+i-1}(Z + W)\gamma'_{l+i-1}(Z + W). \end{aligned}$$

These expansions are summarized in Theorem 9.7. The proof, which is similar to that of Theorem 9.5, is left to the reader.

Theorem 9.7 Let Z be an $m \times m$ symmetric matrix with eigenvalues $x_1 \geq \cdots \geq x_{l-1} > x_l = x_{l+1} = \cdots = x_{l+r-1} > x_{l+r} \geq \cdots \geq x_m$, so that x_l is an eigenvalue with multiplicity r . Suppose that W is an $m \times m$ symmetric matrix, and let $\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_m$ be the eigenvalues of $Z + W$, whereas $\bar{\lambda}_{l,l+r-1} = r^{-1}(\lambda_l + \cdots + \lambda_{l+r-1})$. Denote the eigenprojection of Z corresponding to the repeated eigenvalue x_l by P_l , and denote the total eigenprojection of $Z + W$ corresponding to the collection of eigenvalues $\lambda_l, \dots, \lambda_{l+r-1}$ by Φ_l . Define $Y = (Z - x_l I_m)^+$. Then the third-order Taylor approximations

$$\begin{aligned} \bar{\lambda}_{l,l+r-1} &\approx x_l + a_1 + a_2 + a_3, \\ \Phi_l &\approx P_l + B_1 + B_2 + B_3 \end{aligned}$$

have

$$\begin{aligned}
 a_1 &= \frac{1}{r} \text{tr}(WP_l), \\
 a_2 &= -\frac{1}{r} \text{tr}(WYW P_l), \\
 a_3 &= \frac{1}{r} \{ \text{tr}(YWYW P_l W) - \text{tr}(Y^2 W P_l W P_l W) \}, \\
 B_1 &= -YWP_l - P_l WY, \\
 B_2 &= YWP_l WY + YWYW P_l - Y^2 W P_l W P_l + P_l WY WY \\
 &\quad - P_l W P_l WY^2 - P_l WY^2 W P_l, \\
 B_3 &= Y^2 W P_l WY W P_l + P_l WY W P_l WY^2 + Y^2 W P_l W P_l WY \\
 &\quad + YW P_l W P_l WY^2 + Y^2 WY W P_l W P_l + P_l W P_l WY WY^2 \\
 &\quad + YWY^2 W P_l W P_l + P_l W P_l WY^2 WY - Y^3 W P_l W P_l W P_l \\
 &\quad - P_l W P_l W P_l WY^3 - YWYW P_l WY - YW P_l WY WY \\
 &\quad - YWY WY W P_l - P_l WY WY WY + YW P_l WY^2 W P_l \\
 &\quad + P_l WY^2 W P_l WY + P_l WY^2 WY W P_l + P_l WY WY^2 W P_l \\
 &\quad - P_l WY^3 W P_l W P_l - P_l W P_l WY^3 W P_l.
 \end{aligned}$$

Example 9.4 Suppose that Ω is an $m \times m$ covariance matrix and its smallest eigenvalue, λ , has multiplicity r . Let S be the sample covariance matrix defined in Section 1.13, and define $A = S - \Omega$ so that $S = \Omega + A$. We can use Theorem 9.7 to obtain approximations of functions of S in terms of A . For instance, as an illustration, we will consider the first-order approximation of $\hat{P}(S - \hat{\lambda} I_m) \hat{P}$, where $\hat{\lambda}$ is the average of the r smallest eigenvalues of S and \hat{P} is the total eigenprojection of S corresponding to these r eigenvalues. Now from Theorem 9.7, we have the approximations $\hat{\lambda} \approx \lambda + r^{-1} \text{tr}(AP)$ and $\hat{P} \approx P + B_1$, where P is the eigenprojection of Ω corresponding to its smallest eigenvalue and the formula for B_1 can be obtained from Theorem 9.7. Ignoring second-order terms in A , these approximations yield the approximation

$$\begin{aligned}
 \hat{P}(S - \hat{\lambda} I_m) \hat{P} &\approx (P + B_1)((\Omega + A) - \{\lambda + r^{-1} \text{tr}(AP)\} I_m)(P + B_1) \\
 &= (P + B_1)((\Omega - \lambda I_m) + \{A - r^{-1} \text{tr}(AP) I_m\})(P + B_1) \\
 &= P(\Omega - \lambda I_m)P + B_1(\Omega - \lambda I_m)P + P(\Omega - \lambda I_m)B_1 \\
 &\quad + P\{A - r^{-1} \text{tr}(AP) I_m\}P \\
 &= P\{A - r^{-1} \text{tr}(AP) I_m\}P,
 \end{aligned}$$

where the last equality has used the fact that $P(\Omega - \lambda I_m) = (\Omega - \lambda I_m)P = (0)$.

9.7 MAXIMA AND MINIMA

One important application of derivatives involves finding the maxima or minima of a function. A function f has a local maximum at an $n \times 1$ point \mathbf{a} if for some $\delta > 0$, $f(\mathbf{a}) \geq f(\mathbf{a} + \mathbf{x})$ whenever $\mathbf{x}'\mathbf{x} < \delta$. This function has an absolute maximum at \mathbf{a} if $f(\mathbf{a}) \geq f(\mathbf{x})$ for all \mathbf{x} for which f is defined. Similar definitions hold for a local minimum and an absolute minimum; in fact, if f has a local minimum at a point \mathbf{a} , then $-f$ has a local maximum at \mathbf{a} , and if f has an absolute minimum at \mathbf{a} , then $-f$ has an absolute maximum at \mathbf{a} . For this reason, we will at times confine our discussion to only the case of a maximum. In this section and the next section, we state some results that are helpful in finding local maxima and minima. For proofs of these results, see Khuri (2003) or Magnus and Neudecker (1999). Our first result gives a necessary condition for a function f to have a local maximum at \mathbf{a} .

Theorem 9.8 Suppose the function $f(\mathbf{x})$ is defined for all $n \times 1$ vectors $\mathbf{x} \in S$, where S is some subset of R^n . Let \mathbf{a} be an interior point of S ; that is, a $\delta > 0$ exists, such that $\mathbf{a} + \mathbf{u} \in S$ for all $\mathbf{u}'\mathbf{u} < \delta$. If f has a local maximum at \mathbf{a} and f is differentiable at \mathbf{a} , then

$$\frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) = \mathbf{0}'. \quad (9.29)$$

Any point \mathbf{a} satisfying (9.29) is called a stationary point of f . Although Theorem 9.8 indicates that any point at which a local maximum or local minimum occurs must be a stationary point, the converse does not hold. A stationary point that does not correspond to a local maximum or a local minimum is called a saddle point. Theorem 9.9 is helpful in determining whether a particular stationary point is a local maximum or minimum in those situations in which the function f is twice differentiable.

Theorem 9.9 Suppose the function $f(\mathbf{x})$ is defined for all $n \times 1$ vectors $\mathbf{x} \in S$, where S is some subset of R^n . Suppose also that f is twice differentiable at the interior point \mathbf{a} of S . If \mathbf{a} is a stationary point of f and H_f is the Hessian matrix of f at \mathbf{a} , then

- (a) f has a local minimum at \mathbf{a} if H_f is positive definite,
- (b) f has a local maximum at \mathbf{a} if H_f is negative definite,
- (c) f has a saddle point at \mathbf{a} if H_f is nonsingular but not positive definite or negative definite,
- (d) f may have a local minimum, a local maximum, or a saddle point at \mathbf{a} if H_f is singular.

Example 9.5 On several occasions, we have discussed the problem of finding a least squares solution $\hat{\beta}$ to the inconsistent system of equations

$$\mathbf{y} = \mathbf{X}\beta,$$

where \mathbf{y} is an $N \times 1$ vector of constants, X is an $N \times (k + 1)$ matrix of constants, and β is a $(k + 1) \times 1$ vector of variables. A solution was obtained in Chapter 2 by using the geometrical properties of least squares regression, whereas in Chapter 6 we utilized the results developed on least squares generalized inverses. In this example, we will show how the methods of this section may help obtain a solution. We will assume that $\text{rank}(X) = k + 1$; that is, the matrix X has full column rank. Recall that a least squares solution $\hat{\beta}$ is any vector that minimizes the sum of squared errors given by

$$f(\hat{\beta}) = (\mathbf{y} - X\hat{\beta})'(\mathbf{y} - X\hat{\beta}).$$

The differential of $f(\hat{\beta})$ is

$$\begin{aligned} df &= \{d(\mathbf{y} - X\hat{\beta})'\}(\mathbf{y} - X\hat{\beta}) + (\mathbf{y} - X\hat{\beta})'d(\mathbf{y} - X\hat{\beta}) \\ &= -(d\hat{\beta})'X'(\mathbf{y} - X\hat{\beta}) - (\mathbf{y} - X\hat{\beta})'Xd\hat{\beta} \\ &= -2(\mathbf{y} - X\hat{\beta})'Xd\hat{\beta}, \end{aligned}$$

so that

$$\frac{\partial}{\partial \hat{\beta}} f(\hat{\beta}) = -2(\mathbf{y} - X\hat{\beta})'X.$$

Thus, on setting this first derivative equal to $\mathbf{0}'$ and rearranging, we find that the stationary values are given by the solutions $\hat{\beta}$ to the system of equations

$$X'X\hat{\beta} = X'\mathbf{y}. \quad (9.30)$$

Since X has full column rank, $X'X$ is nonsingular, and so the unique solution to (9.30) is

$$\hat{\beta} = (X'X)^{-1}X'\mathbf{y}. \quad (9.31)$$

To verify that this solution minimizes the sum of squared errors, we need to obtain the Hessian matrix H_f . The second differential of $f(\hat{\beta})$ is given by

$$\begin{aligned} d^2f &= d(df) = -d\{2(\mathbf{y} - X\hat{\beta})'Xd\hat{\beta}\} \\ &= -2\{d(\mathbf{y} - X\hat{\beta})'\}Xd\hat{\beta} \\ &= 2(d\hat{\beta})'X'Xd\hat{\beta}, \end{aligned}$$

so that

$$H_f = 2X'X.$$

Since this matrix is positive definite, it follows from Theorem 9.9 that the solution given in (9.31) minimizes $f(\hat{\beta})$.

Example 9.6 One of the most popular ways of obtaining estimators of unknown parameters is by a method known as maximum likelihood estimation. If we have a random sample of vectors $\mathbf{x}_1, \dots, \mathbf{x}_n$ from a population having density function $f(\mathbf{x}; \boldsymbol{\theta})$, where $\boldsymbol{\theta}$ is a vector of parameters, then the likelihood function of $\boldsymbol{\theta}$ is defined to be the joint density function of $\mathbf{x}_1, \dots, \mathbf{x}_n$ viewed as a function of $\boldsymbol{\theta}$; that is, this likelihood function is given by

$$L(\boldsymbol{\theta}) = \prod_{i=1}^n f(\mathbf{x}_i; \boldsymbol{\theta}).$$

The method of maximum likelihood estimates $\boldsymbol{\theta}$ by the vector $\hat{\boldsymbol{\theta}}$ that maximizes $L(\boldsymbol{\theta})$. In this example, we will use this method to obtain estimates of $\boldsymbol{\mu}$ and Ω when our sample is coming from the normal distribution, $N_m(\boldsymbol{\mu}, \Omega)$. Thus, $\boldsymbol{\mu}$ is an $m \times 1$ vector, Ω is an $m \times m$ positive definite matrix, and the required density function $f(\mathbf{x}; \boldsymbol{\mu}, \Omega)$ is given in (1.13). In deriving the estimates $\hat{\boldsymbol{\mu}}$ and $\hat{\Omega}$, we will find it a little bit easier to maximize the function $\log(L(\boldsymbol{\mu}, \Omega))$, which is, of course, maximized at the same solution as is $L(\boldsymbol{\mu}, \Omega)$. After omitting terms from $\log(L(\boldsymbol{\mu}, \Omega))$ that do not involve $\boldsymbol{\mu}$ or Ω , we find that we must maximize the function

$$g(\boldsymbol{\mu}, \Omega) = -\frac{1}{2}n \log |\Omega| - \frac{1}{2}\text{tr}(\Omega^{-1}U),$$

where

$$U = \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})'.$$

The first differential of g is given by

$$\begin{aligned} dg &= -\frac{1}{2}n d(\log |\Omega|) - \frac{1}{2}\text{tr}\{(d\Omega^{-1})U\} - \frac{1}{2}\text{tr}(\Omega^{-1}dU) \\ &= -\frac{1}{2}n \text{tr}(\Omega^{-1}d\Omega) + \frac{1}{2}\text{tr}\{\Omega^{-1}(d\Omega)\Omega^{-1}U\} \\ &\quad + \frac{1}{2}\text{tr}\left(\Omega^{-1}\left\{(d\boldsymbol{\mu})\sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})' + \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})d\boldsymbol{\mu}'\right\}\right) \\ &= \frac{1}{2}\text{tr}\{(d\Omega)\Omega^{-1}(U - n\Omega)\Omega^{-1}\} \\ &\quad + \frac{1}{2}\text{tr}(\Omega^{-1}\{n(d\boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu})d\boldsymbol{\mu}'\}) \\ &= \frac{1}{2}\text{tr}\{(d\Omega)\Omega^{-1}(U - n\Omega)\Omega^{-1}\} + n(\bar{\mathbf{x}} - \boldsymbol{\mu})'\Omega^{-1}d\boldsymbol{\mu} \\ &= \frac{1}{2}\text{vec}(d\Omega)'(\Omega^{-1} \otimes \Omega^{-1}) \text{vec}(U - n\Omega) + n(\bar{\mathbf{x}} - \boldsymbol{\mu})'\Omega^{-1}d\boldsymbol{\mu}, \end{aligned}$$

where the second equality applied Corollary 9.1.1 and Theorem 9.2, and the fifth equality applied Theorem 8.12. Since Ω is symmetric, $\text{vec}(d\Omega) = d \text{vec}(\Omega) =$

$D_m d \mathbf{v}(\Omega)$, and so the differential may be re-expressed as

$$dg = \frac{1}{2} \{d \mathbf{v}(\Omega)\}' D_m' (\Omega^{-1} \otimes \Omega^{-1}) \text{vec}(U - n\Omega) + n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \Omega^{-1} d\boldsymbol{\mu}, \quad (9.32)$$

and thus,

$$\begin{aligned} \frac{\partial}{\partial \boldsymbol{\mu}'} g &= n(\bar{\mathbf{x}} - \boldsymbol{\mu})' \Omega^{-1}, \\ \frac{\partial}{\partial \mathbf{v}(\Omega)'} g &= \frac{1}{2} \{ \text{vec}(U - n\Omega) \}' (\Omega^{-1} \otimes \Omega^{-1}) D_m'. \end{aligned}$$

On equating these first derivatives to null vectors, we obtain the equations

$$\begin{aligned} n\Omega^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}) &= \mathbf{0}, \\ D_m'(\Omega^{-1} \otimes \Omega^{-1}) \text{vec}(U - n\Omega) &= \mathbf{0}. \end{aligned}$$

From the first of these two equations, we obtain the solution for $\boldsymbol{\mu}$, $\hat{\boldsymbol{\mu}} = \bar{\mathbf{x}}$, whereas the second can be rewritten as

$$D_m'(\Omega^{-1} \otimes \Omega^{-1}) D_m \mathbf{v}(U - n\Omega) = \mathbf{0},$$

because the symmetry of $(U - n\Omega)$ implies that $\text{vec}(U - n\Omega) = D_m \mathbf{v}(U - n\Omega)$. Premultiplying this equation by $D_m^+(\Omega \otimes \Omega) D_m^{+'}$ and using Theorem 8.35, we find that

$$\mathbf{v}(U - n\Omega) = \mathbf{0}.$$

Since $(U - n\Omega)$ is symmetric, this implies that $(U - n\Omega) = (0)$, and so the solution for Ω is $\hat{\Omega} = n^{-1}U$. All that remains is to show that the solution $(\hat{\boldsymbol{\mu}}, \hat{\Omega})$ yields a maximum. By differentiating (9.32), we find that

$$\begin{aligned} d^2g &= \frac{1}{2} \{d \mathbf{v}(\Omega)\}' D_m' \{d(\Omega^{-1} \otimes \Omega^{-1})\} \text{vec}(U - n\Omega) \\ &\quad + \frac{1}{2} \{d \mathbf{v}(\Omega)\}' D_m' (\Omega^{-1} \otimes \Omega^{-1}) \text{vec}(dU - nd\Omega) \\ &\quad - n(d\boldsymbol{\mu})' \Omega^{-1} d\boldsymbol{\mu} + n(\bar{\mathbf{x}} - \boldsymbol{\mu})' (d\Omega^{-1}) d\boldsymbol{\mu}. \end{aligned}$$

Evaluating this at $\boldsymbol{\mu} = \bar{\mathbf{x}}$ and $\Omega = n^{-1}U$, we find that the first and the fourth terms on the right-hand side of the equation above vanish. In addition, note that

$$dU = n(d\boldsymbol{\mu})(\bar{\mathbf{x}} - \boldsymbol{\mu})' + n(\bar{\mathbf{x}} - \boldsymbol{\mu}) d\boldsymbol{\mu}'$$

also vanishes when evaluated at $\boldsymbol{\mu} = \bar{\mathbf{x}}$. Thus, at $\boldsymbol{\mu} = \bar{\mathbf{x}}$ and $\Omega = n^{-1}U$,

$$\begin{aligned} d^2g &= -\frac{n}{2} \{d \mathbf{v}(\Omega)\}' D_m' (\Omega^{-1} \otimes \Omega^{-1}) D_m d \mathbf{v}(\Omega) - n(d\boldsymbol{\mu})' \Omega^{-1} d\boldsymbol{\mu} \\ &= [d\boldsymbol{\mu}' \quad \{d \mathbf{v}(\Omega)\}'] H_g \begin{bmatrix} d\boldsymbol{\mu} \\ d \mathbf{v}(\Omega) \end{bmatrix}, \end{aligned}$$

where

$$H_g = \begin{bmatrix} -n\Omega^{-1} & (0) \\ (0) & -\frac{n}{2}D'_m(\Omega^{-1} \otimes \Omega^{-1})D_m \end{bmatrix}.$$

Clearly, H_g is negative definite because Ω^{-1} and $D'_m(\Omega^{-1} \otimes \Omega^{-1})D_m$ are positive definite matrices, which then establishes that the solution $(\hat{\mu}, \hat{\Omega}) = (\bar{x}, n^{-1}U)$ yields a maximum.

9.8 CONVEX AND CONCAVE FUNCTIONS

In Section 2.11, we discussed convex sets. Here we will extend the concept of convexity to functions and obtain some special results that apply to this class of functions.

Definition 9.1 Let $f(x)$ be a real-valued function defined for all $x \in S$, where S is a convex subset of R^m . Then $f(x)$ is a convex function on S , if

$$f(cx_1 + (1-c)x_2) \leq cf(x_1) + (1-c)f(x_2) \quad (9.33)$$

for all $x_1 \neq x_2$, $x_1 \in S$, $x_2 \in S$, and $0 < c < 1$. The function $f(x)$ is a strictly convex function if inequality (9.33) is always a strict inequality. If $-f(x)$ is a (strictly) convex function, then $f(x)$ is said to be a (strictly) concave function.

The implication of Definition 9.1 is that the line between any two points of a convex function must be above the function. Figure 9.1 illustrates a convex function when x is a scalar variable.

If $f(x)$ is a convex function, then it is easily verified that the set defined by

$$T = \{z = (x', y)' : x \in S, y \geq f(x)\}$$

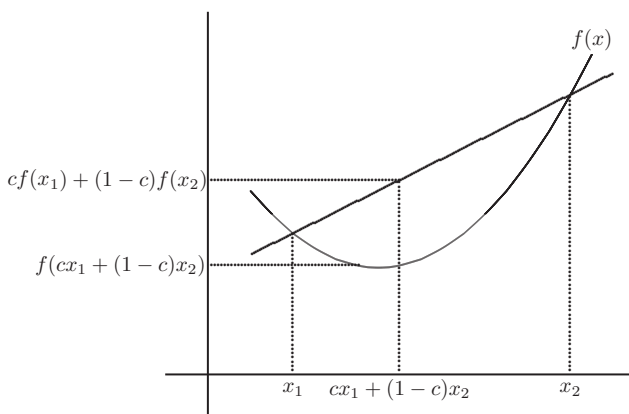


Figure 9.1 A convex function $f(x)$ of a scalar variable x

is a convex subset of R^{m+1} . For instance, if $m = 1$, then T will be a convex subset of R^2 . In this case, for any $a \in S$, the point $(a, f(a))$ will be a boundary point of the set T . Now from the supporting hyperplane theorem, Theorem 2.32, we know that there is a line passing through the point $(a, f(a))$, such that the function $f(x)$ is never below this line. Since this line passes through the point $(a, f(a))$, it can be written in the form $g(x) = f(a) + t(x - a)$, where t is the slope of the line, and thus, for all $x \in S$, we have

$$f(x) \geq f(a) + t(x - a). \quad (9.34)$$

The generalization of this result to arbitrary m is given below.

Theorem 9.10 Let $f(x)$ be a real-valued convex function defined for all $x \in S$, where S is a convex subset of R^m . Then, corresponding to each interior point $a \in S$, an $m \times 1$ vector t exists, such that

$$f(x) \geq f(a) + t'(x - a) \quad (9.35)$$

for all $x \in S$.

Proof. For any $a \in S$, the point $z_* = (a', f(a))'$ is a boundary point of the convex set T defined above, and so it follows from Theorem 2.32 that an $(m + 1) \times 1$ vector $b = (b'_1, b_{m+1})' \neq 0$ exists, for which $b'z \geq b'z_*$ for all $z \in T$. Clearly, for any $z = (x', y) \in T$, we can arbitrarily increase the value of y and get another point in T . For this reason, we see that b_{m+1} cannot be negative because if it was, we could make $b'z$ arbitrarily small and, in particular, less than $b'z_*$. Thus, b_{m+1} is either positive or 0. Now for any $x \in S$, $(x', f(x))' \in T$, and so for this choice of z in the inequality $b'z \geq b'z_*$, we get

$$b'_1 x + b_{m+1} f(x) \geq b'_1 a + b_{m+1} f(a).$$

If b_{m+1} is positive, then the inequality above may be rearranged to the form given in (9.35) with $t = -b_{m+1}^{-1} b'_1$. If, on the other hand, $b_{m+1} = 0$, then $b'z \geq b'z_*$ reduces to

$$b'_1 x \geq b'_1 a,$$

which implies that a is a boundary point of S . Thus, the proof is complete. \square

If f is a differentiable function, then the hyperplane given on the right-hand side of (9.35) will be given by the tangent hyperplane to $f(x)$ at $x = a$.

Theorem 9.11 Let $f(x)$ be a real-valued convex function defined for all $x \in S$, where S is an open convex subset of R^m . If $f(x)$ is differentiable and $a \in S$, then

$$f(x) \geq f(a) + \left(\frac{\partial}{\partial a'} f(a) \right) (x - a)$$

for all $x \in S$.

Proof. Suppose that $\mathbf{x} \in S$ and $\mathbf{a} \in S$, and let $\mathbf{y} = \mathbf{x} - \mathbf{a}$, so that $\mathbf{x} = \mathbf{a} + \mathbf{y}$. Since S is convex, the point

$$c\mathbf{x} + (1 - c)\mathbf{a} = c(\mathbf{a} + \mathbf{y}) + (1 - c)\mathbf{a} = \mathbf{a} + c\mathbf{y}$$

is in S for $0 \leq c \leq 1$. Thus, because of the convexity of f , we have

$$f(\mathbf{a} + c\mathbf{y}) \leq cf(\mathbf{a} + \mathbf{y}) + (1 - c)f(\mathbf{a}) = f(\mathbf{a}) + c\{f(\mathbf{a} + \mathbf{y}) - f(\mathbf{a})\},$$

or equivalently,

$$f(\mathbf{a} + \mathbf{y}) \geq f(\mathbf{a}) + c^{-1}\{f(\mathbf{a} + c\mathbf{y}) - f(\mathbf{a})\}. \quad (9.36)$$

Now because f is differentiable, we also have the Taylor formula

$$f(\mathbf{a} + c\mathbf{y}) = f(\mathbf{a}) + \left(\frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) \right) c\mathbf{y} + r_1(c\mathbf{y}, \mathbf{a}), \quad (9.37)$$

where the remainder satisfies $\lim c^{-1}r_1(c\mathbf{y}, \mathbf{a}) = 0$ as $c \rightarrow 0$. Using (9.37) in (9.36), we get

$$f(\mathbf{a} + \mathbf{y}) \geq f(\mathbf{a}) + \left(\frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) \right) \mathbf{y} + c^{-1}r_1(c\mathbf{y}, \mathbf{a}),$$

and so the result follows by letting $c \rightarrow 0$. □

We can easily use Theorem 9.11 to show that a stationary point of a convex function will actually be an absolute minimum. Equivalently, a stationary point of a concave function will be an absolute maximum of that function.

Theorem 9.12 Let $f(\mathbf{x})$ be a real-valued convex function defined for all $\mathbf{x} \in S$, where S is an open convex subset of R^m . If $f(\mathbf{x})$ is differentiable and $\mathbf{a} \in S$ is a stationary point of f , then f has an absolute minimum at \mathbf{a} .

Proof. If \mathbf{a} is a stationary point of f , then

$$\frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) = \mathbf{0}'.$$

Using this in the inequality of Theorem 9.11, we get $f(\mathbf{x}) \geq f(\mathbf{a})$ for all $\mathbf{x} \in S$, and so the result follows. □

We can use the inequality given in (9.35) to prove a useful inequality involving the moments of a random vector \mathbf{y} . This inequality is known as Jensen's inequality. However, before we can prove this result, we will need Theorem 9.13.

Theorem 9.13 Suppose that S is a convex subset of R^m and \mathbf{y} is an $m \times 1$ random vector with finite first moments. If $P(\mathbf{y} \in S) = 1$, then $E(\mathbf{y}) \in S$.

Proof. We will prove the result by induction. Clearly, the result holds if $m = 1$, because in this case, S is an interval, and it is easily shown that a random variable y satisfying $P(a \leq y \leq b) = 1$ for some constants a and b will have $a \leq E(y) \leq b$. Now assuming that the result holds for dimension $m - 1$, we will show that it must then hold for m . Define the convex set $S_* = \{\mathbf{x} : \mathbf{x} = \mathbf{u} - E(\mathbf{y}), \mathbf{u} \in S\}$, so that the proof will be complete if we show that $\mathbf{0} \in S_*$. Now if $\mathbf{0} \notin S_*$, it follows from Theorem 2.32 that an $m \times 1$ vector $\mathbf{a} \neq \mathbf{0}$ exists, such that $\mathbf{a}'\mathbf{x} \geq 0$ for all $\mathbf{x} \in S_*$. Consequently, because $P(\mathbf{y} \in S) = P(\mathbf{w} \in S_*) = 1$, where the random vector $\mathbf{w} = \mathbf{y} - E(\mathbf{y})$, we have $\mathbf{a}'\mathbf{w} \geq 0$ with probability 1; yet $E(\mathbf{a}'\mathbf{w}) = 0$, which is possible only if $\mathbf{a}'\mathbf{w} = 0$, in which case, \mathbf{w} is on the hyperplane defined by $\{\mathbf{x} : \mathbf{a}'\mathbf{x} = 0\}$, with probability one. However, because $P(\mathbf{w} \in S_*) = 1$ as well, we must have $P(\mathbf{w} \in S_0) = 1$, where $S_0 = S_* \cap \{\mathbf{x} : \mathbf{a}'\mathbf{x} = 0\}$. Now it follows from Theorem 2.28 that S_0 is a convex set, and it is contained within an $(m - 1)$ -dimensional vector space because $\{\mathbf{x} : \mathbf{a}'\mathbf{x} = 0\}$ is an $(m - 1)$ -dimensional vector space. Thus, because our result holds for $(m - 1)$ -dimensional spaces, we must have $E(\mathbf{w}) = \mathbf{0} \in S_0$. This leads to the contradiction $\mathbf{0} \in S_*$, because $S_0 \subseteq S_*$, and so the proof is complete. \square

We now prove Jensen's inequality in Theorem 9.14.

Theorem 9.14 Let $f(x)$ be a real-valued convex function defined for all $x \in S$, where S is a convex subset of R^m . If \mathbf{y} is an $m \times 1$ random vector with finite first moments and satisfying $P(\mathbf{y} \in S) = 1$, then

$$E(f(\mathbf{y})) \geq f(E(\mathbf{y})).$$

Proof. Theorem 9.13 guarantees that $E(\mathbf{y}) \in S$. We first prove the result for $m = 1$. If $E(y)$ is an interior point of S , the result follows by taking the expected value of both sides of (9.34) when $x = y$ and $a = E(y)$. Since S is an interval when $m = 1$, $E(y)$ can be a boundary point of S only if S is closed and $P(y = c) = 1$, where c is an endpoint of the interval. In this case, the result is trivial because the terms on the two sides of the inequality above are equal. We will complete the proof by showing that if the result holds for $m - 1$, then it must hold for m . If the $m \times 1$ vector $E(\mathbf{y})$ is an interior point of S , the result follows by taking the expected value of both sides of (9.35) with $\mathbf{x} = \mathbf{y}$ and $\mathbf{a} = E(\mathbf{y})$. If $E(\mathbf{y})$ is a boundary point of S , then we know from the supporting hyperplane theorem that an $m \times 1$ unit vector \mathbf{b} exists, such that $w = \mathbf{b}'\mathbf{y} \geq \mathbf{b}'E(\mathbf{y}) = \mu$ with probability one. However, because we also have $E(w) = \mathbf{b}'E(\mathbf{y}) = \mu$, it follows that $\mathbf{b}'\mathbf{y} = \mu$ with probability one. Let P be any $m \times m$ orthogonal matrix with its last column given by \mathbf{b} , so that the vector $\mathbf{u} = P'\mathbf{y}$ has the form $\mathbf{u} = (\mathbf{u}_1, \mu)'$, where \mathbf{u}_1 is an $(m - 1) \times 1$ vector. Define the function $g(\mathbf{u}_1)$ as

$$g(\mathbf{u}_1) = f\left(P \begin{bmatrix} \mathbf{u}_1 \\ \mu \end{bmatrix}\right) = f(\mathbf{y}),$$

for all $\mathbf{u}_1 \in S_* = \{\mathbf{x} : \mathbf{x} = P'_1 \mathbf{y}, \mathbf{y} \in S\}$, where P_1 is the matrix obtained from P by deleting its last column. The convexity of S_* and g follow from the convexity of S and f , and so, because \mathbf{u}_1 is $(m-1) \times 1$, our result applies to $g(\mathbf{u}_1)$. Thus, we have

$$\begin{aligned} E(f(\mathbf{y})) &= E(g(\mathbf{u}_1)) \geq g(E(\mathbf{u}_1)) \\ &= f\left(P \begin{bmatrix} E(\mathbf{u}_1) \\ \mu \end{bmatrix}\right) = f(E(\mathbf{y})), \end{aligned}$$

and so the proof is complete. \square

9.9 THE METHOD OF LAGRANGE MULTIPLIERS

In some situations, we may need to find a local maximum of a function $f(\mathbf{x})$, where f is defined for all $\mathbf{x} \in S$, whereas the desired maximum is over all \mathbf{x} in T , a subset of S . The method of Lagrange multipliers is useful in those situations in which the set T can be expressed in terms of a number of equality constraints; that is, functions g_1, \dots, g_m exist, such that

$$T = \{\mathbf{x} : \mathbf{x} \in R^n, \mathbf{g}(\mathbf{x}) = \mathbf{0}\},$$

where $\mathbf{g}(\mathbf{x})$ is the $m \times 1$ function given by $(g_1(\mathbf{x}), \dots, g_m(\mathbf{x}))'$.

The method of Lagrange multipliers involves the maximization of the Lagrange function

$$L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \boldsymbol{\lambda}' \mathbf{g}(\mathbf{x}),$$

where the components of the $m \times 1$ vector $\boldsymbol{\lambda}$, $\lambda_1, \dots, \lambda_m$, are called the Lagrange multipliers. The stationary values of $L(\mathbf{x}, \boldsymbol{\lambda})$ are the solutions $(\mathbf{x}, \boldsymbol{\lambda})$ satisfying

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}'} L(\mathbf{x}, \boldsymbol{\lambda}) &= \frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) - \boldsymbol{\lambda}' \left(\frac{\partial}{\partial \mathbf{x}'} \mathbf{g}(\mathbf{x}) \right) = \mathbf{0}', \\ \frac{\partial}{\partial \boldsymbol{\lambda}'} L(\mathbf{x}, \boldsymbol{\lambda}) &= -\mathbf{g}(\mathbf{x})' = \mathbf{0}'. \end{aligned} \tag{9.38}$$

The second equation above is simply the equality constraints,

$$\mathbf{g}(\mathbf{x}) = \mathbf{0} \tag{9.39}$$

that determine the set T . Under certain conditions, the local maximum of the function $f(\mathbf{x})$, subject to $\mathbf{x} \in T$, will be given by a vector \mathbf{x} that, for some $\boldsymbol{\lambda}$, satisfies equations (9.38) and (9.39). We will present a procedure for determining whether a particular solution vector \mathbf{x} is a local maximum. This procedure is based on the following result, a proof of which can be found in Magnus and Neudecker (1999).

Theorem 9.15 Suppose the function $f(\mathbf{x})$ is defined for all $n \times 1$ vectors $\mathbf{x} \in S$, where S is some subset of R^n , and $\mathbf{g}(\mathbf{x})$ is an $m \times 1$ vector function defined for all $\mathbf{x} \in S$, where $m < n$. Let \mathbf{a} be an interior point of S , and suppose that the following conditions hold:

- (a) f and \mathbf{g} are twice differentiable at \mathbf{a} .
- (b) The first derivative of \mathbf{g} at \mathbf{a} , $(\partial/\partial \mathbf{a}')\mathbf{g}(\mathbf{a})$, has full rank m .
- (c) $\mathbf{g}(\mathbf{a}) = \mathbf{0}$.
- (d) $(\partial/\partial \mathbf{a}')L(\mathbf{a}, \boldsymbol{\lambda}) = \mathbf{0}'$, where $L(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) - \boldsymbol{\lambda}'\mathbf{g}(\mathbf{x})$ and $\boldsymbol{\lambda}$ is $m \times 1$.

Let H_f and H_{g_i} be the Hessian matrices of the functions $f(\mathbf{x})$ and $g_i(\mathbf{x})$ evaluated at $\mathbf{x} = \mathbf{a}$, and define

$$A = H_f - \sum_{i=1}^m \lambda_i H_{g_i},$$

$$B = \frac{\partial}{\partial \mathbf{a}'} \mathbf{g}(\mathbf{a}).$$

Then $f(\mathbf{x})$ has a local maximum at $\mathbf{x} = \mathbf{a}$, subject to $\mathbf{g}(\mathbf{x}) = \mathbf{0}$, if

$$\mathbf{x}' A \mathbf{x} < 0,$$

for all $\mathbf{x} \neq \mathbf{0}$ for which $B\mathbf{x} = \mathbf{0}$.

A similar result holds for a local minimum with the inequality $\mathbf{x}' A \mathbf{x} > 0$ replacing $\mathbf{x}' A \mathbf{x} < 0$. Theorem 9.16 provides a method for determining whether $\mathbf{x}' A \mathbf{x} < 0$ or $\mathbf{x}' A \mathbf{x} > 0$ holds for all $\mathbf{x} \neq \mathbf{0}$ satisfying $B\mathbf{x} = \mathbf{0}$. Again, a proof can be found in Magnus and Neudecker (1999).

Theorem 9.16 Let A be an $n \times n$ symmetric matrix and B be an $m \times n$ matrix. For $r = 1, \dots, n$, let A_{rr} be the $r \times r$ matrix obtained by deleting the last $n - r$ rows and columns of A , and let B_r be the $m \times r$ matrix obtained by deleting the last $n - r$ columns of B . For $r = 1, \dots, n$, define the $(m + r) \times (m + r)$ matrix Δ_r as

$$\Delta_r = \begin{bmatrix} (0) & B_r \\ B_r' & A_{rr} \end{bmatrix}.$$

Then, if B_m is nonsingular, $\mathbf{x}' A \mathbf{x} > 0$ holds for all $\mathbf{x} \neq \mathbf{0}$ satisfying $B\mathbf{x} = \mathbf{0}$ if and only if

$$(-1)^m |\Delta_r| > 0,$$

for $r = m + 1, \dots, n$, and $\mathbf{x}' A \mathbf{x} < 0$ holds for all $\mathbf{x} \neq \mathbf{0}$ satisfying $B\mathbf{x} = \mathbf{0}$ if and only if

$$(-1)^r |\Delta_r| > 0,$$

for $r = m + 1, \dots, n$.

Example 9.7 We will find solutions $\mathbf{x} = (x_1, x_2, x_3)'$, which maximize and minimize the function

$$f(\mathbf{x}) = x_1 + x_2 + x_3,$$

subject to the constraints

$$x_1^2 + x_2^2 = 1, \quad (9.40)$$

$$x_3 - x_1 - x_2 = 1. \quad (9.41)$$

Setting the first derivative of the Lagrange function

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

with respect to \mathbf{x} , equal to $\mathbf{0}'$, we obtain the equations

$$1 - 2\lambda_1 x_1 + \lambda_2 = 0,$$

$$1 - 2\lambda_1 x_2 + \lambda_2 = 0,$$

$$1 - \lambda_2 = 0.$$

The third equation gives $\lambda_2 = 1$, and when this is substituted in the first two equations, we find that we must have

$$x_1 = x_2 = \frac{1}{\lambda_1}.$$

Using this equation in (9.40), we find that $\lambda_1 = \pm\sqrt{2}$, and so we have the stationary points

$$(x_1, x_2, x_3) = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}}, 1 + \sqrt{2} \right) \quad \text{when } \lambda_1 = \sqrt{2},$$

$$(x_1, x_2, x_3) = \left(-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}, 1 - \sqrt{2} \right) \quad \text{when } \lambda_1 = -\sqrt{2}.$$

To determine whether either of these solutions yields a maximum or minimum, we use Theorem 9.15 and Theorem 9.16. Thus, because $m = 2$ and $n = 3$, we only need the determinant of the matrix

$$\Delta_3 = \begin{bmatrix} 0 & 0 & 2x_1 & 2x_2 & 0 \\ 0 & 0 & -1 & -1 & 1 \\ 2x_1 & -1 & -2\lambda_1 & 0 & 0 \\ 2x_2 & -1 & 0 & -2\lambda_1 & 0 \\ 0 & 1 & 0 & 0 & 0 \end{bmatrix}.$$

By using the cofactor expansion formula for a determinant, it is fairly straightforward to show that

$$|\Delta_3| = -8\lambda_1(x_1^2 + x_2^2).$$

Thus, when $(x_1, x_2, x_3, \lambda_1, \lambda_2) = (1/\sqrt{2}, 1/\sqrt{2}, 1 + \sqrt{2}, \sqrt{2}, 1)$, we have

$$(-1)^r |\Delta_r| = (-1)^3 |\Delta_3| = 8\sqrt{2} > 0,$$

and so the solution $(x_1, x_2, x_3) = (1/\sqrt{2}, 1/\sqrt{2}, 1 + \sqrt{2})$ yields a constrained maximum. On the other hand, when $(x_1, x_2, x_3, \lambda_1, \lambda_2) = (-1/\sqrt{2}, -1/\sqrt{2}, 1 - \sqrt{2}, -\sqrt{2}, 1)$,

$$(-1)^m |\Delta_r| = (-1)^2 |\Delta_3| = 8\sqrt{2} > 0,$$

so the solution $(x_1, x_2, x_3) = (-1/\sqrt{2}, -1/\sqrt{2}, 1 - \sqrt{2})$ yields a constrained minimum.

In some situations, in the process of obtaining the stationary values of $L(\mathbf{x}, \boldsymbol{\lambda})$, it becomes apparent which solution yields a maximum and which solution yields a minimum. Thus, in this case, there will be no need to compute the Δ_r matrices.

Example 9.8 Let A be an $m \times m$ symmetric matrix and \mathbf{x} be an $m \times 1$ nonnull vector. We saw in Section 3.6 that

$$\frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}} \tag{9.42}$$

has a maximum value of $\lambda_1(A)$ and a minimum value of $\lambda_m(A)$, where $\lambda_1(A) \geq \dots \geq \lambda_m(A)$ are the eigenvalues of A . We will prove this result again, this time using Lagrange's method. Note that because $\mathbf{z} = (\mathbf{x}' \mathbf{x})^{-1/2} \mathbf{x}$ is a unit vector, it follows that maximizing or minimizing (9.42) over all $\mathbf{x} \neq \mathbf{0}$ is equivalent to maximizing or minimizing the function

$$f(\mathbf{z}) = \mathbf{z}' A \mathbf{z},$$

subject to the constraint

$$\mathbf{z}' \mathbf{z} = 1. \tag{9.43}$$

Thus, the Lagrange function is

$$L(\mathbf{z}, \lambda) = \mathbf{z}' A \mathbf{z} - \lambda(\mathbf{z}' \mathbf{z} - 1).$$

Setting its first derivative, with respect to \mathbf{z}' , equal to $\mathbf{0}'$, we obtain the equation

$$2A\mathbf{z} - 2\lambda\mathbf{z} = \mathbf{0},$$

or, equivalently,

$$A\mathbf{z} = \lambda\mathbf{z}, \tag{9.44}$$

which is the eigenvalue-eigenvector equation for A . Thus, the Lagrange multiplier λ is an eigenvalue of A . Further, premultiplying (9.44) by z' and using (9.43), we find that

$$\lambda = z'Az;$$

that is, if $(z', \lambda)'$ is a stationary point of $L(z, \lambda)$, then $\lambda = z'Az$ must be an eigenvalue of A . Consequently, the maximum value of $z'Az$, subject to $z'z = 1$, is $\lambda_1(A)$, which is attained when z is equal to any unit eigenvector corresponding to $\lambda_1(A)$. Similarly, the minimum value of $z'Az$, subject to $z'z = 1$, is $\lambda_m(A)$, which is attained at any unit eigenvector associated with $\lambda_m(A)$.

Example 9.9 Let $D = \text{diag}(d_1, \dots, d_m)$, where $d_1 > \dots > d_m > 0$ and let A be an $m \times m$ positive definite matrix. Consider the function $f(X) = \text{tr}(XAX'D)$, where X is an $m \times m$ orthogonal matrix. We wish to maximize and minimize $f(X)$ over all choices for X . Since X must satisfy $XX' - I_m = (0)$, we use the Lagrange function

$$L(X, \Lambda) = \text{tr}(XAX'D) + \text{tr}\{\Lambda(XX' - I_m)\},$$

where Λ is a symmetric matrix of Lagrange multipliers. The differential of $L(X, \Lambda)$ is

$$\begin{aligned} dL &= \text{tr}\{(dX)AX'D\} + \text{tr}\{XA(dX)'D\} + \text{tr}\{d(\Lambda)(XX' - I_m)\} \\ &\quad + \text{tr}\{\Lambda(dX)X'\} + \text{tr}\{\Lambda X(dX)'\} \\ &= 2 \text{vec}(DXA + \Lambda X)'d \text{vec}(X) + \text{vec}(XX' - I_m)'d \text{vec}(\Lambda), \end{aligned}$$

and so the stationary values of $L(X, \Lambda)$ occur at the solutions to

$$DXA + \Lambda X = (0), \quad XX' = I_m.$$

Thus, $\Lambda = -DXAX'$, and since Λ is symmetric, we have $DXAX' = XAX'D$, or

$$DY = YD, \tag{9.45}$$

where $Y = XAX'$. Examining the (i, j) th term on each side of (9.45), we see that $d_i y_{ij} = y_{ij} d_j$. Since $d_i \neq d_j$ when $i \neq j$, it follows that Y is a diagonal matrix. Consequently, the stationary values of $f(X)$ subject to $XX' = I_m$ are given by the set of values

$$\sum_{i=1}^m d_{j_i} \lambda_i(A),$$

where $\lambda_1(A) \geq \dots \geq \lambda_m(A)$ are the eigenvalues of A and (j_1, \dots, j_m) is a permutation of $(1, \dots, m)$, the set being formed over all such permutations. It readily then follows that

$$\max_{XX'=I_m} \text{tr}(XAX'D) = \sum_{i=1}^m d_i \lambda_i(A),$$

and

$$\min_{X X' = I_m} \text{tr}(X A X' D) = \sum_{i=1}^m d_{m-i+1} \lambda_i(A).$$

In Example 9.10, we obtain the best quadratic unbiased estimator of σ^2 in the ordinary least squares regression model.

Example 9.10 Consider the multiple regression model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 I)$. A quadratic estimator of σ^2 is any estimator $\hat{\sigma}^2$ that takes the form $\hat{\sigma}^2 = \mathbf{y}' A \mathbf{y}$, where A is a symmetric matrix of constants. We wish to find the choice of A that minimizes $\text{var}(\hat{\sigma}^2)$ over all choices of A for which $\hat{\sigma}^2$ is unbiased. Now because $E(\boldsymbol{\epsilon}) = \mathbf{0}$ and $E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}') = \sigma^2 I$, we have

$$\begin{aligned} E(\mathbf{y}' A \mathbf{y}) &= E\{(X\boldsymbol{\beta} + \boldsymbol{\epsilon})' A (X\boldsymbol{\beta} + \boldsymbol{\epsilon})\} \\ &= E\{\boldsymbol{\beta}' X' A X \boldsymbol{\beta} + 2\boldsymbol{\beta}' X' A \boldsymbol{\epsilon} + \boldsymbol{\epsilon}' A \boldsymbol{\epsilon}\} \\ &= \boldsymbol{\beta}' X' A X \boldsymbol{\beta} + \text{tr}\{A E(\boldsymbol{\epsilon}\boldsymbol{\epsilon}')\} \\ &= \boldsymbol{\beta}' X' A X \boldsymbol{\beta} + \sigma^2 \text{tr}(A), \end{aligned}$$

and so $\hat{\sigma}^2 = \mathbf{y}' A \mathbf{y}$ is unbiased regardless of the value of $\boldsymbol{\beta}$ only if

$$X' A X = (0) \quad (9.46)$$

and

$$\text{tr}(A) = 1. \quad (9.47)$$

Using the fact that the components of $\boldsymbol{\epsilon}$ are independently distributed and the first four moments of each component are 0, 1, 0, 3, it is easily verified that

$$\text{var}(\mathbf{y}' A \mathbf{y}) = 2\sigma^4 \text{tr}(A^2) + 4\sigma^2 \boldsymbol{\beta}' X' A^2 X \boldsymbol{\beta}.$$

Thus, the required Lagrange function is

$$L(A, \lambda, \Lambda) = 2\sigma^4 \text{tr}(A^2) + 4\sigma^2 \boldsymbol{\beta}' X' A^2 X \boldsymbol{\beta} - \text{tr}(\Lambda X' A X) - \lambda \{\text{tr}(A) - 1\},$$

where the Lagrange multipliers are given by λ and the components of the matrix Λ , which is symmetric because $X' A X$ is symmetric. Differentiation with respect to A yields

$$\begin{aligned} dL &= 2\sigma^4 \text{tr}\{(dA)A + A dA\} + 4\sigma^2 \boldsymbol{\beta}' X' \{(dA)A + A dA\} X \boldsymbol{\beta} \\ &\quad - \text{tr}\{\Lambda X' (dA) X\} - \lambda \text{tr}(dA) \\ &= \text{tr}\{4\sigma^4 A + 4\sigma^2 (A X \boldsymbol{\beta} \boldsymbol{\beta}' X' + X \boldsymbol{\beta} \boldsymbol{\beta}' X' A) - X \Lambda X' - \lambda I_N\} dA. \end{aligned}$$

Thus, we must use

$$4\sigma^4 A + 4\sigma^2 (A X \boldsymbol{\beta} \boldsymbol{\beta}' X' + X \boldsymbol{\beta} \boldsymbol{\beta}' X' A) - X \Lambda X' - \lambda I_N = (0) \quad (9.48)$$

along with (9.46) and (9.47) to solve for A . Premultiplying and postmultiplying (9.48) by XX^+ and using (9.46) and the fact that $X^+ = (X'X)^+X'$, we find that

$$X\Lambda X' = -\lambda XX^+.$$

Substituting this result back into (9.48), we get

$$A = \frac{1}{4}\sigma^{-4}\lambda(I_N - XX^+) - \sigma^{-2}H, \quad (9.49)$$

where $H = A\gamma\gamma' + \gamma\gamma'A$ and $\gamma = X\beta$. Putting (9.49) back into (9.48) and simplifying, we obtain

$$H = -\sigma^{-2}(H\gamma\gamma' + \gamma\gamma'H). \quad (9.50)$$

By postmultiplying (9.50) by γ , we find that γ must be an eigenvector of H , which in light of (9.50), can be true only if H is of the form $H = c\gamma\gamma'$ for some scalar c . Further, when we put $H = c\gamma\gamma'$ in (9.50), we find that we must have $c = 0$; thus, $H = (0)$. In addition, if we take the trace of both sides of (9.49) and use (9.47), we see that

$$\lambda = \frac{4\sigma^4}{\text{tr}(I_N - XX^+)} = \frac{4\sigma^4}{N - r},$$

where r is the rank of X . Consequently, we have shown that (9.49) simplifies to

$$A = (N - r)^{-1}(I_N - XX^+), \quad (9.51)$$

so that $\hat{\sigma}^2 = \mathbf{y}'A\mathbf{y} = \text{SSE}/(N - r)$ is the familiar residual variance estimate. We can easily show that (9.51) yields an absolute minimum by writing an arbitrary symmetric matrix satisfying (9.46) and (9.47), as $A_* = A + B$, where B must then satisfy $\text{tr}(B) = 0$ and $X'BX = (0)$. Then, because $\text{tr}(AB) = 0$ and $AX = (0)$, we have

$$\begin{aligned} \text{var}(\mathbf{y}'A_*\mathbf{y}) &= 2\sigma^4\text{tr}(A_*^2) + 4\sigma^2\beta'X'A_*^2X\beta \\ &= 2\sigma^4\{\text{tr}(A^2) + \text{tr}(B^2) + 2\text{tr}(AB)\} + 4\sigma^2\beta'X' \\ &\quad \times (A^2 + B^2 + AB + BA)X\beta \\ &= 2\sigma^4\{\text{tr}(A^2) + \text{tr}(B^2)\} + 4\sigma^2\beta'X'B^2X\beta \\ &\geq 2\sigma^4\text{tr}(A^2) = \text{var}(\mathbf{y}'A\mathbf{y}). \end{aligned}$$

PROBLEMS

9.1 Consider the natural log function, $f(x) = \log(x)$.

- (a) Obtain the k th-order Taylor formula for $f(1 + u)$ in powers of u .
- (b) Use the formula in part (a) with $k = 5$ to approximate $\log(1.1)$.

9.2 Suppose the function f of the 2×1 vector \mathbf{x} is given by

$$f(\mathbf{x}) = \frac{(x_2 - 1)^2}{(x_1 + 1)^3}.$$

Give the second-order Taylor formula for $f(\mathbf{0} + \mathbf{u})$ in powers of u_1 and u_2 .

9.3 Suppose the 2×1 function \mathbf{f} of the 3×1 vector \mathbf{x} is given by

$$\mathbf{f}(\mathbf{x}) = \begin{bmatrix} x_1^2 + x_2^2 + x_3^2 \\ 2x_1 - x_2 - x_3 \end{bmatrix}$$

and the 2×1 function \mathbf{g} of the 2×1 vector \mathbf{z} is given by

$$\mathbf{g}(\mathbf{z}) = \begin{bmatrix} z_2/z_1 \\ z_1 z_2 \end{bmatrix}.$$

Use the chain rule to compute

$$\frac{\partial}{\partial \mathbf{x}'} \mathbf{y}(\mathbf{x}),$$

where $\mathbf{y}(\mathbf{x})$ is the composite function defined by $\mathbf{y}(\mathbf{x}) = \mathbf{g}(\mathbf{f}(\mathbf{x}))$.

9.4 Let A be an $m \times m$ symmetric matrix of constants and \mathbf{x} be an $m \times 1$ vector of variables. Find the differential and first derivative of $f(\mathbf{x}) = \exp(-\frac{1}{2}\mathbf{x}'A\mathbf{x})$.

9.5 Let A and B be $m \times m$ symmetric matrices of constants and \mathbf{x} be an $m \times 1$ vector of variables. Find the differential and first derivative of the function

$$f(\mathbf{x}) = \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'B\mathbf{x}}.$$

9.6 Let \mathbf{x} be an $m \times 1$ vector of variables. Find the differential and derivative of $\mathbf{x}\mathbf{x}'$.

9.7 Let A and B be $m \times n$ matrices of constants and X be an $n \times m$ matrix of variables. Find the differential and derivative of

(a) $\text{tr}(AX)$,

(b) $\text{tr}(AXBX)$.

9.8 Let X be an $m \times m$ nonsingular matrix, A be an $m \times m$ matrix of constants, and \mathbf{a} be an $m \times 1$ vector of constants. Find the differential and derivative of

(a) $|X^2|$,

(b) $\text{tr}(AX^{-1})$,

(c) $\mathbf{a}'X^{-1}\mathbf{a}$.

9.9 Let A be an $m \times m$ positive definite matrix of constants and X be an $m \times n$ matrix of variables such that $X'AX$ is nonsingular. Find the differential and first derivative of $\log |X'AX|$.

9.10 Let X be an $m \times n$ matrix with $\text{rank}(X) = n$. Show that

$$\frac{\partial}{\partial \text{vec}(X)'} |X'X| = 2|X'X| (\text{vec}\{X(X'X)^{-1}\})'.$$

9.11 Let X be an $m \times m$ matrix and n be a positive integer. Show that

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(X^n) = \sum_{i=1}^n \{(X^{n-i})' \otimes X^{i-1}\}.$$

9.12 Let A and B be $n \times m$ and $m \times n$ matrices of constants, respectively. If X is an $m \times m$ nonsingular matrix, find the derivatives of

- (a) $\text{vec}(AXB)$,
- (b) $\text{vec}(AX^{-1}B)$.

9.13 Show that if X is an $m \times m$ nonsingular matrix and $X_{\#}$ is its adjoint matrix, then

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(X_{\#}) = |X| \{ \text{vec}(X^{-1}) \text{vec}(X^{-1'})' - (X^{-1'} \otimes X^{-1}) \}.$$

9.14 Prove Corollary 9.1.1.

9.15 Let X be an $m \times m$ symmetric matrix of variables. For each of the following functions, find the Jacobian matrix

$$\frac{\partial}{\partial \text{v}(X)'} \text{vec}(F).$$

- (a) $F(X) = AXA'$, where A is an $m \times m$ matrix of constants.
- (b) $F(X) = XBX$, where B is an $m \times m$ symmetric matrix of constants.

9.16 Suppose X is an $m \times n$ matrix of rank n . Find the differential and first derivative of $I_m - X(X'X)^{-1}X'$.

9.17 Let X be an $m \times n$ matrix. Show that

- (a) if $F(X) = X \otimes X$,

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(F) = (I_n \otimes K_{nm} \otimes I_m) \{ I_{mn} \otimes \text{vec}(X) + \text{vec}(X) \otimes I_{mn} \},$$

- (b) if $F(X) = X \odot X$,

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(F) = 2D_{\text{vec}(X)}.$$

9.18 Show that the Hessian matrix H_f is given by

- (a) $H_f = 2I_{mn}$ if $f(X) = \text{tr}(X'X)$, where X is an $m \times n$ matrix,
- (b) $H_f = 2K_{mm}$ if $f(X) = \text{tr}(X^2)$, where X is an $m \times m$ matrix,

- (c) $H_f = -K_{mm}(X^{-1'} \otimes X^{-1})$ if $f(X) = \log |X|$, where X is an $m \times m$ nonsingular matrix.

9.19 Let X be an $m \times m$ nonsingular matrix. Show that

$$d^n X^{-1} = (-1)^n n! (X^{-1} dX)^n X^{-1}.$$

9.20 Let X be an $m \times m$ matrix having correlation structure; that is, X is a symmetric matrix of variables, except that each of its diagonal elements is equal to one. Show that if X is nonsingular, then

$$\frac{\partial}{\partial \tilde{v}(X)'} \tilde{v}(X^{-1}) = -2 \tilde{L}_m N_m (X^{-1} \otimes X^{-1}) N_m \tilde{L}_m'.$$

9.21 Suppose that Y is an $m \times m$ symmetric matrix and ϵ is a scalar, such that $(I_m + \epsilon Y)^{-1}$ exists. Let $(I_m + \epsilon Y)^{-1/2}$ be the symmetric square root of $(I_m + \epsilon Y)^{-1}$, so that

$$(I_m + \epsilon Y)^{-1} = (I_m + \epsilon Y)^{-1/2} (I_m + \epsilon Y)^{-1/2}.$$

Using perturbation methods, show that

$$(I_m + \epsilon Y)^{-1/2} = I_m + \sum_{i=1}^{\infty} \epsilon^i B_i,$$

where

$$B_1 = -\frac{1}{2}Y, \quad B_2 = \frac{3}{8}Y^2, \quad B_3 = -\frac{5}{16}Y^3, \quad B_4 = \frac{35}{128}Y^4.$$

9.22 Let X be an $m \times n$ full column rank matrix, so that $X^+ = (X'X)^{-1}X'$. Let Y be an $m \times n$ matrix and ϵ be a scalar, such that $X + \epsilon Y$ is also full column rank. Show that

$$(X + \epsilon Y)^+ = X^+ + \sum_{i=1}^{\infty} \epsilon^i B_i,$$

where

$$B_1 = (X'X)^{-1}Y'(I_m - XX^+) - X^+YX^+.$$

9.23 Let S be an $m \times m$ sample covariance matrix, and suppose that Ω , the corresponding population covariance matrix, has each of its diagonal elements equal to one. Define A to be the difference between these two matrices; that is, $A = S - \Omega$, so that $S = \Omega + A$. Note that the population correlation matrix is also Ω , whereas the sample correlation matrix is given by $R = D_S^{-1/2} S D_S^{-1/2}$, where $D_S^{-1/2} = \text{diag}(s_{11}^{-1/2}, \dots, s_{mm}^{-1/2})$. Show that the approximation

$R = \Omega + C_1 + C_2 + C_3$, accurate up through third-order terms in the elements of A , is given by

$$\begin{aligned} C_1 &= A - \frac{1}{2}(\Omega D_A + D_A \Omega), \\ C_2 &= \frac{3}{8}(D_A^2 \Omega + \Omega D_A^2) + \frac{1}{4}D_A \Omega D_A - \frac{1}{2}(AD_A + D_A A), \\ C_3 &= \frac{3}{8}(D_A^2 A + AD_A^2) + \frac{1}{4}D_A AD_A - \frac{3}{16}(D_A^2 \Omega D_A + D_A \Omega D_A^2) \\ &\quad - \frac{5}{16}(D_A^3 \Omega + \Omega D_A^3), \end{aligned}$$

where $D_A = \text{diag}(a_{11}, \dots, a_{mm})$.

9.24 Derive the results given in Theorem 9.7. First obtain expressions for B_1 , B_2 , and B_3 by using the equations $(Z + W)\Phi_l = \Phi_l(Z + W)$, $\Phi_l^2 = \Phi_l$, and $\Phi_l' = \Phi_l$. Then obtain expressions for a_1 , a_2 , and a_3 by using the fact that $\bar{\lambda}_{l,l+r-1} = r^{-1} \text{tr}\{(Z + W)\Phi_l\}$.

9.25 Let $X = \text{diag}(x_1, \dots, x_m)$, where $x_1 \geq \dots \geq x_m$, and suppose that the l th diagonal element is distinct, so that $x_l \neq x_i$ if $i \neq l$. Let $\lambda_1 \geq \dots \geq \lambda_m$ and $\gamma_1, \dots, \gamma_m$ be the eigenvalues and corresponding eigenvectors of $(I_m + V)^{-1}(X + U)$, where U and V are $m \times m$ symmetric matrices; that is, for each i ,

$$(X + U)\gamma_i = \lambda_i(I_m + V)\gamma_i.$$

The purpose of this exercise is to obtain the first-order approximations $\lambda_l = x_l + a_1$ and $\gamma_l = ce_l + b_1$, where e_l is the l th column of I_m . Higher order approximations can be found in Sugiura (1976). These approximations can be determined by using the eigenvalue-eigenvector equation just given along with the appropriate scale constraint on γ_l .

(a) Show that $a_1 = u_{ll} - x_l v_{ll}$.

(b) Show that if $c = 1$ and $\gamma_l' \gamma_l = 1$, then

$$b_{l1} = 0, \quad b_{i1} = -\frac{u_{il} - x_l v_{il}}{x_i - x_l} \quad \text{for all } i \neq l,$$

where b_{i1} is the i th component of the vector b_1 .

(c) Show that if $c = 1$ and $\gamma_l'(I_m + V)\gamma_l = 1$, then

$$b_{l1} = -\frac{1}{2}v_{ll}, \quad b_{i1} = -\frac{u_{il} - x_l v_{il}}{x_i - x_l} \quad \text{for all } i \neq l.$$

(d) Show that if $c = x_l^{1/2}$ and $\gamma_l' \gamma_l = \lambda_l$, then

$$b_{l1} = \frac{u_{ll} - x_l v_{ll}}{2x_l^{1/2}}, \quad b_{i1} = -\frac{x_l^{1/2}(u_{il} - x_l v_{il})}{x_i - x_l} \quad \text{for all } i \neq l.$$

- 9.26** Let Ω and S be as defined in Example 9.4 with the smallest eigenvalue of Ω being $\lambda = 1$, and consider the quantity

$$U = \frac{r^{-1} \sum_{i=m-r+1}^m \lambda_i^2(S)}{\{r^{-1} \sum_{i=m-r+1}^m \lambda_i(S)\}^2} - 1,$$

where $\lambda_1(S) \geq \cdots \geq \lambda_m(S)$ are the eigenvalues of S . Note that

$$\sum_{i=m-r+1}^m \lambda_i^2(S) = \text{tr}(S^2 \hat{P}),$$

where \hat{P} is the total eigenprojection of S corresponding to its r smallest eigenvalues. Show that if we let $A = S - \Omega$, so that $S = \Omega + A$, then the second-order approximation formula for U is given by

$$U \approx r^{-1}(\text{tr}(APAP) - r^{-1}\{\text{tr}(AP)\}^2).$$

- 9.27** Consider the function f of the 2×1 vector \mathbf{x} given by

$$f(\mathbf{x}) = 2x_1^3 + x_2^3 - 6x_1 - 27x_2.$$

- (a) Determine the stationary points of f .
 - (b) Identify each of the points in part (a) as a maximum, minimum, or saddle point.
- 9.28** For each of the following functions, determine any local maxima or minima:
- (a) $x_1^2 + \frac{1}{2}x_2^2 - 2x_1x_2 + x_1 - 2x_2 + 1$.
 - (b) $x_1^3 + \frac{3}{2}x_1^2 + x_2^2 - 6x_1 - 2x_2$.
 - (c) $x_2^3 + 2x_1^2 + x_3^2 + 2x_1x_3 - 3x_2 - x_3$.
- 9.29** Let \mathbf{a} be an $m \times 1$ vector and B be an $m \times m$ symmetric matrix, each containing constants. Let \mathbf{x} be an $m \times 1$ vector of variables.

- (a) Show that the function

$$f(\mathbf{x}) = \mathbf{x}'B\mathbf{x} + \mathbf{a}'\mathbf{x}$$

has stationary solutions given by

$$\mathbf{x} = -\frac{1}{2}B^+\mathbf{a} + (I_m - B^+B)\mathbf{y},$$

where \mathbf{y} is an arbitrary $m \times 1$ vector.

- (b) Show that if B is nonsingular, then only one stationary solution exists. When will this solution yield a maximum or a minimum?

9.30 If the Hessian matrix H_f of a function f is singular at a stationary point \mathbf{x} , then we must take a closer look at the behavior of this function in the neighborhood of the point \mathbf{x} to determine whether the point is a maximum, minimum, or a saddle point. For each of the functions below, show that $\mathbf{0}$ is a stationary point and the Hessian matrix is singular at $\mathbf{0}$. In each case, determine whether $\mathbf{0}$ yields a maximum, minimum, or a saddle point.

(a) $x_1^4 + x_2^4$.

(b) $x_1^2 x_2^2 - x_1^4 - x_2^4$.

(c) $x_1^3 - x_2^3$.

9.31 Suppose that we have independent random samples from each of k multivariate normal distributions with the i th distribution being $N_m(\boldsymbol{\mu}_i, \Omega)$. Thus, these distributions have possibly different mean vectors but identical covariance matrices. If the i th sample is denoted by $\mathbf{x}_{i1}, \dots, \mathbf{x}_{in_i}$, show that the maximum likelihood estimators of $\boldsymbol{\mu}_i$ and Ω are given by

$$\hat{\boldsymbol{\mu}}_i = \bar{\mathbf{x}}_i = \sum_{j=1}^{n_i} \frac{\mathbf{x}_{ij}}{n_i}, \quad \hat{\Omega} = \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)(\mathbf{x}_{ij} - \bar{\mathbf{x}}_i)'}{n},$$

where $n = n_1 + \dots + n_k$.

9.32 Consider the multiple regression model,

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where \mathbf{y} is $N \times 1$, X is $N \times m$, $\boldsymbol{\beta}$ is $m \times 1$, and $\boldsymbol{\epsilon}$ is $N \times 1$. Suppose that $\text{rank}(X) = m$ and $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 I_N)$, so that $\mathbf{y} \sim N_N(X\boldsymbol{\beta}, \sigma^2 I_N)$. Find the maximum likelihood estimates of $\boldsymbol{\beta}$ and σ^2 .

9.33 Let $f(\mathbf{x})$ be a real-valued convex function defined for all $\mathbf{x} \in S$, where S is a convex subset of R^m . Show that the set $T = \{\mathbf{z} = (\mathbf{x}', y)' : \mathbf{x} \in S, y \geq f(\mathbf{x})\}$ is convex.

9.34 Suppose that $f(\mathbf{x})$ and $g(\mathbf{x})$ are convex functions both defined on the convex set $S \subseteq R^m$. Show that the function $af(\mathbf{x}) + bg(\mathbf{x})$ is convex if a and b are nonnegative scalars.

9.35 Prove the converse of Theorem 9.11; that is, show that if $f(\mathbf{x})$ is defined and differentiable on the open convex set S and

$$f(\mathbf{x}) \geq f(\mathbf{a}) + \left(\frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) \right) (\mathbf{x} - \mathbf{a})$$

for all $\mathbf{x} \in S$ and $\mathbf{a} \in S$, then $f(\mathbf{x})$ is a convex function.

9.36 Let $f(\mathbf{x})$ be a real-valued function defined for all $\mathbf{x} \in S$, where S is an open convex subset of R^m , and suppose that $f(\mathbf{x})$ is a twice differentiable function on S . Show that $f(\mathbf{x})$ is a convex function if and only if the Hessian matrix H_f is nonnegative definite at each $\mathbf{x} \in S$.

9.37 Let \mathbf{x} be a 2×1 vector and consider the function $f(\mathbf{x}) = x_1^c x_2^{1-c}$ for all $\mathbf{x} \in S$, where $0 < c < 1$ and $S = \{\mathbf{x} : x_1 > 0, x_2 > 0\}$.

(a) Use the previous exercise to show that $f(\mathbf{x})$ is a concave function.

(b) Show that if \mathbf{y} is a 2×1 random vector with finite first moments and satisfying $P(\mathbf{y} \in S) = 1$, then

$$E(y_1^\alpha y_2^{1-\alpha}) \leq \{E(y_1)\}^\alpha \{E(y_2)\}^{1-\alpha}$$

if $0 < \alpha < 1$.

9.38 Let \mathbf{x} be a 3×1 vector, and define the function

$$f(\mathbf{x}) = x_1 + x_2 - x_3.$$

Find the maximum and minimum of $f(\mathbf{x})$ subject to the constraint $\mathbf{x}'\mathbf{x} = 1$.

9.39 Find the shortest distance from the origin to a point on the surface given by

$$x_1^2 + x_2^2 + x_3^2 + 4x_1 - 6x_3 = 2.$$

9.40 Let A be an $m \times m$ positive definite matrix and \mathbf{x} be an $m \times 1$ vector. Find the maximum and minimum of the function

$$f(\mathbf{x}) = \mathbf{x}'\mathbf{x},$$

subject to the constraint $\mathbf{x}'A\mathbf{x} = 1$.

9.41 Find the maximum and minimum of the function

$$f(\mathbf{x}) = x_1(x_2 + x_3),$$

subject to the constraints $x_1^2 + x_2^2 = 1$ and $x_1x_3 + x_2 = 2$.

9.42 For an $m \times 1$ vector \mathbf{x} with positive components, maximize the function

$$f(\mathbf{x}) = x_1x_2 \cdots x_m,$$

subject to the constraint $x_1 + x_2 + \cdots + x_m = a$, where a is some positive number. Use this to establish the inequality

$$(x_1x_2 \cdots x_m)^{1/m} \leq \frac{1}{m}(x_1 + \cdots + x_m)$$

for all positive real numbers x_1, \dots, x_m .

9.43 Let A and B be $m \times m$ matrices, with A being nonnegative definite and B being positive definite. Following the approach of Example 9.8, apply the Lagrange method to find the maximum and minimum values of

$$f(\mathbf{x}) = \frac{\mathbf{x}'A\mathbf{x}}{\mathbf{x}'B\mathbf{x}},$$

over all $\mathbf{x} \neq \mathbf{0}$.

- 9.44** Let \mathbf{a} be an $m \times 1$ nonnull vector and B be an $m \times m$ positive definite matrix. Using the results of Problem 9.43, show that for $\mathbf{x} \neq \mathbf{0}$,

$$f(\mathbf{x}) = \frac{(\mathbf{a}'\mathbf{x})^2}{\mathbf{x}'B\mathbf{x}}$$

has a maximum value of

$$\mathbf{a}'B^{-1}\mathbf{a}.$$

We can use this result to obtain the union-intersection test (see Example 3.16) of the multivariate hypothesis $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$ against $H_1 : \boldsymbol{\mu} \neq \boldsymbol{\mu}_0$, where $\boldsymbol{\mu}$ represents the $m \times 1$ mean vector of a population and $\boldsymbol{\mu}_0$ is an $m \times 1$ vector of constants. Let $\bar{\mathbf{x}}$ and S denote the sample mean vector and sample covariance matrix computed from a sample of size n from this population. Show that if we base the union-intersection procedure on the univariate t statistic

$$t = \frac{(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)'}{s/\sqrt{n}}$$

for testing $H_0 : \boldsymbol{\mu} = \boldsymbol{\mu}_0$, then the union-intersection test can be based on $T^2 = n(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)'S^{-1}(\bar{\mathbf{x}} - \boldsymbol{\mu}_0)$.

- 9.45** Suppose that x_1, \dots, x_n are independent and identically distributed random variables with mean μ and variance σ^2 . Consider a linear estimator of μ , which is any estimator of the form $\hat{\mu} = \sum a_i x_i$, where a_1, \dots, a_n are constants.
- (a) For what values of a_1, \dots, a_n will $\hat{\mu}$ be an unbiased estimator of μ ?
- (b) Use the method of Lagrange multipliers to show that the sample mean \bar{x} is the best linear unbiased estimator of μ ; that is, \bar{x} has the smallest variance among all unbiased linear estimators of μ .
- 9.46** A random process involves n independent trials, where each trial can result in one of k distinct outcomes. Let p_i denote the probability that a trial results in outcome i and note that then $p_1 + \dots + p_k = 1$. Define the random variables, x_1, \dots, x_k , where x_i counts the number of times that outcome i occurs in the n trials. Then the random vector $\mathbf{x} = (x_1, \dots, x_k)'$ has the multinomial distribution with probability function given by

$$P(x_1 = n_1, \dots, x_k = n_k) = \frac{n!}{n_1! \dots n_k!} p_1^{n_1} \dots p_k^{n_k},$$

where n_1, \dots, n_k are nonnegative integers satisfying $n_1 + \dots + n_k = n$. Find the maximum likelihood estimate of $\mathbf{p} = (p_1, \dots, p_k)'$.

- 9.47** Suppose that the $m \times m$ positive definite covariance matrix Ω is partitioned in the form

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}' & \Omega_{22} \end{bmatrix},$$

where Ω_{11} is $m_1 \times m_1$, Ω_{22} is $m_2 \times m_2$, and $m_1 + m_2 = m$. Suppose also that the $m \times 1$ random vector \mathbf{x} has covariance matrix Ω and is partitioned as

$\mathbf{x} = (\mathbf{x}'_1, \mathbf{x}'_2)'$, where \mathbf{x}_1 is $m_1 \times 1$ and \mathbf{x}_2 is $m_2 \times 1$. If the $m_1 \times 1$ vector \mathbf{a} and $m_2 \times 1$ vector \mathbf{b} are vectors of constants, then the square of the correlation between the random variables $u = \mathbf{a}'\mathbf{x}_1$ and $v = \mathbf{b}'\mathbf{x}_2$ is given by

$$f(\mathbf{a}, \mathbf{b}) = \frac{(\mathbf{a}'\Omega_{12}\mathbf{b})^2}{\mathbf{a}'\Omega_{11}\mathbf{a}\mathbf{b}'\Omega_{22}\mathbf{b}}.$$

Show that the maximum value of $f(\mathbf{a}, \mathbf{b})$, that is, the maximum squared correlation between u and v , subject to the constraints

$$\mathbf{a}'\Omega_{11}\mathbf{a} = 1, \quad \mathbf{b}'\Omega_{22}\mathbf{b} = 1,$$

is the largest eigenvalue of $\Omega_{11}^{-1}\Omega_{12}\Omega_{22}^{-1}\Omega'_{12}$ or, equivalently, the largest eigenvalue of $\Omega_{22}^{-1}\Omega'_{12}\Omega_{11}^{-1}\Omega_{12}$. What are the vectors \mathbf{a} and \mathbf{b} that yield this maximum?

9.48 Consider the multiple regression model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, where $E(\boldsymbol{\epsilon}) = \mathbf{0}$, $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 I_N$, and the $N \times (k+1)$ matrix X has rank $k+1$. A linear estimator of $\boldsymbol{\beta}$ has the form $\tilde{\boldsymbol{\beta}} = L\mathbf{y}$, where L is a $(k+1) \times N$ matrix of constants.

- (a) What is the condition on L for $\tilde{\boldsymbol{\beta}}$ to be an unbiased estimator of $\boldsymbol{\beta}$?
- (b) Show that the minimum variance unbiased linear estimator of $\boldsymbol{\beta}$ is the least squares estimator $\hat{\boldsymbol{\beta}} = (X'X)^{-1}X'\mathbf{y}$. That is, show that $\text{tr}\{\text{var}(\tilde{\boldsymbol{\beta}})\}$ is minimized by choosing $L = (X'X)^{-1}X'$.

10

INEQUALITIES

10.1 INTRODUCTION

We have already encountered a number of inequalities in this text. For instance, we saw the Cauchy-Schwarz inequality in Chapter 2, the Poincaré separation theorem in Chapter 3, the Hadamard inequality in Chapter 8 and Jensen's inequality in Chapter 9, just to name a few of them. In this chapter, we present some additional important classical inequalities along with some extensions or generalizations of them. We begin the chapter with an introduction to the theory of majorization, a concept which has proven useful in developing a number of inequalities.

10.2 MAJORIZATION

Majorization is a preordering of the vectors in R^m . It can be viewed as an attempt to make precise the notion that the components of one vector \mathbf{x} are less spread out than the components of another vector \mathbf{y} . As an example, we may be comparing the vector $\mathbf{x} = m^{-1}\mathbf{1}_m$ to the vector $\mathbf{y} = \mathbf{e}_{1,m}$. In this section, we develop some of the basic results regarding majorization and a few of its applications to the eigenvalues of symmetric matrices. Additional results and applications can be found in Marshall et al. (2011), Ando (1989, 1994), and Bhatia (1997).

In defining majorization, we will need to refer to the ordered components of vectors. For a vector $\mathbf{x} = (x_1, \dots, x_m)'$, denote its decreasing components by $x_{[1]} \geq \dots \geq x_{[m]}$ and its increasing components by $x_{(1)} \leq \dots \leq x_{(m)}$.

Definition 10.1 A vector \mathbf{x} is said to be majorized by a vector \mathbf{y} , indicated by $\mathbf{x} \prec \mathbf{y}$, and \mathbf{y} is said to majorize \mathbf{x} , indicated by $\mathbf{y} \succ \mathbf{x}$, if

$$\sum_{i=1}^k x_{[i]} \leq \sum_{i=1}^k y_{[i]}, \quad k = 1, \dots, m-1,$$

$$\sum_{i=1}^m x_i = \sum_{i=1}^m y_i,$$

or, equivalently,

$$\sum_{i=1}^k x_{(i)} \geq \sum_{i=1}^k y_{(i)}, \quad k = 1, \dots, m-1,$$

$$\sum_{i=1}^m x_i = \sum_{i=1}^m y_i.$$

Clearly, for any \mathbf{x} , $\mathbf{x} \prec \mathbf{x}$, and for any \mathbf{x} , \mathbf{y} , and \mathbf{z} satisfying $\mathbf{x} \prec \mathbf{y}$ and $\mathbf{y} \prec \mathbf{z}$, then $\mathbf{x} \prec \mathbf{z}$. Also, notice that a majorization relationship between \mathbf{x} and \mathbf{y} is not dependent on the ordering of the components of \mathbf{x} or \mathbf{y} . So, for instance, if P and Q are $m \times m$ permutation matrices, then $\mathbf{x} \prec \mathbf{y}$ if and only if $P\mathbf{x} \prec Q\mathbf{y}$. As a result, for notational convenience, in some situations we will assume without loss of generality that the components of \mathbf{x} and \mathbf{y} have the ordering $x_1 \geq \dots \geq x_m$ and $y_1 \geq \dots \geq y_m$.

As an example, consider the $m \times 1$ vectors $\mathbf{x}_k = k^{-1} \sum_{i=1}^k \mathbf{e}_{i,m}$, $k = 1, \dots, m$. Then $m^{-1}\mathbf{1}_m = \mathbf{x}_m \prec \mathbf{x}_{m-1} \prec \dots \prec \mathbf{x}_1 = \mathbf{e}_{1,m}$, and for any $m \times 1$ vector \mathbf{x} having nonnegative components and satisfying $\mathbf{1}'_m \mathbf{x} = 1$, $\mathbf{x}_m \prec \mathbf{x} \prec \mathbf{x}_1$. However, a majorization relationship does not exist for every pair of vectors. For instance, if $\mathbf{x} = \frac{2}{3}\mathbf{e}_{1,m} + \frac{1}{6}\mathbf{e}_{2,m} + \frac{1}{6}\mathbf{e}_{3,m}$, then \mathbf{x} does not majorize \mathbf{x}_2 nor does \mathbf{x}_2 majorize \mathbf{x} .

There is an important connection between majorization and doubly stochastic matrices. An $m \times m$ matrix P is doubly stochastic if it has nonnegative components and $P\mathbf{1}_m = P'\mathbf{1}_m = \mathbf{1}_m$. We first establish the relation between $\mathbf{x} = P\mathbf{y}$ and \mathbf{y} .

Theorem 10.1 If \mathbf{y} is an $m \times 1$ vector and P is an $m \times m$ doubly stochastic matrix, then $\mathbf{x} = P\mathbf{y}$ is majorized by \mathbf{y} .

Proof. We assume $y_1 \geq \dots \geq y_m$ since, otherwise, we may replace \mathbf{y} and P by $Q\mathbf{y}$ and PQ' , where Q is a permutation matrix for which $Q\mathbf{y}$ has decreasing components. Since $\mathbf{1}'_m P = \mathbf{1}'_m$, it immediately follows that

$$\sum_{i=1}^m x_i = \mathbf{1}'_m \mathbf{x} = \mathbf{1}'_m \mathbf{y} = \sum_{i=1}^m y_i.$$

Note that $x_{[i]} = \sum_{j=1}^m p_{h_i j} y_j$ for some choice of h_i , so $\sum_{i=1}^k x_{[i]} = \sum_{j=1}^m w_{jk} y_j$, where $w_{jk} = \sum_{i=1}^k p_{h_i j} \leq \sum_{i=1}^m p_{ij} = 1$ and $\sum_{j=1}^m w_{jk} = k$. Consequently, we have

$$\begin{aligned}
 \sum_{i=1}^k y_{[i]} - \sum_{i=1}^k x_{[i]} &= \sum_{i=1}^k y_i - \sum_{i=1}^m w_{ik} y_i \\
 &= \sum_{i=1}^k (1 - w_{ik}) y_i + \sum_{i=1}^k w_{ik} y_i - \sum_{i=1}^m w_{ik} y_i \\
 &= \sum_{i=1}^k (1 - w_{ik}) y_i - \sum_{i=k+1}^m w_{ik} y_i \\
 &\geq y_k \left(k - \sum_{i=1}^k w_{ik} \right) - y_{k+1} \sum_{i=k+1}^m w_{ik} \\
 &= (y_k - y_{k+1}) \sum_{i=k+1}^m w_{ik} \geq 0,
 \end{aligned}$$

and so the proof is complete. \square

Our next result will make use of a special linear transformation known as a T -transform. It utilizes the matrix

$$T = \lambda I_m + (1 - \lambda)Q,$$

where $0 \leq \lambda \leq 1$ and Q is a permutation matrix that interchanges two components. It is easily observed that T is doubly stochastic since Q is doubly stochastic. Theorem 10.1 indicates that the relationship $\mathbf{x} = P\mathbf{y}$, where P is doubly stochastic, is a sufficient condition for $\mathbf{x} \prec \mathbf{y}$. Theorem 10.2 tells us it is also a necessary condition.

Theorem 10.2 If \mathbf{x} is majorized by \mathbf{y} , then there exists a doubly stochastic matrix P such that $\mathbf{x} = P\mathbf{y}$.

Proof. We assume without loss of generality that $x_1 \geq \cdots \geq x_m$ and $y_1 \geq \cdots \geq y_m$, and let r denote the number of nonzero differences $y_i - x_i$. If $r = 0$, the result is trivial since $\mathbf{x} = I_m \mathbf{y}$. We will prove the result by induction, assuming it holds when $r < h$ and proving that it holds when $r = h$. Thus, for the remainder of the proof, we have $\mathbf{x} \prec \mathbf{y}$ with exactly $h > 0$ nonzero differences $y_i - x_i$. We will find a T -transform matrix T such that $\mathbf{x} \prec \mathbf{z} = T\mathbf{y}$ and the number of nonzero differences $z_i - x_i$ is less than h . Then the assumed result will yield $\mathbf{x} = P\mathbf{z} = PT\mathbf{y}$ for some doubly stochastic matrix P and the proof will be complete since PT is also doubly stochastic. Let j be the largest index for which $x_j < y_j$ and k be the smallest

index greater than j for which $x_k > y_k$. Such a j must exist since $r = h > 0$ and \mathbf{y} majorizes \mathbf{x} , while a corresponding k is guaranteed by the fact that $x_i > y_i$ for the largest index i satisfying $x_i \neq y_i$. As a result, we have

$$y_j > x_j \geq x_k > y_k, \quad (10.1)$$

and

$$x_i = y_i, \quad i = j + 1, \dots, k - 1. \quad (10.2)$$

Let $\lambda = 1 - \delta/(y_j - y_k)$, where $\delta = \min(y_j - x_j, x_k - y_k)$, so that $0 < \lambda < 1$ follows from (10.1). Then using $T = \lambda I_m + (1 - \lambda)Q$, where Q is the permutation matrix obtained from the identity matrix by interchanging the j th and k th rows, we find that

$$\begin{aligned} \mathbf{z} = T\mathbf{y} &= (y_1, \dots, y_{j-1}, \lambda y_j + (1 - \lambda)y_k, y_{j+1}, \dots, y_{k-1}, \\ &\quad \lambda y_k + (1 - \lambda)y_j, y_{k+1}, \dots, y_m)' \\ &= (y_1, \dots, y_{j-1}, y_j - \delta, y_{j+1}, \dots, y_{k-1}, y_k + \delta, y_{k+1}, \dots, y_m)'. \end{aligned}$$

We will now show that $\mathbf{x} \prec \mathbf{z}$. First note that since $x_j \leq z_j = y_j - \delta$ and $x_k \geq z_k = y_k + \delta$, we have

$$\begin{aligned} z_{j-1} = y_{j-1} &\geq y_j > z_j \geq x_j \geq x_{j+1} = y_{j+1} = z_{j+1}, \\ z_{k-1} = y_{k-1} &= x_{k-1} \geq x_k \geq z_k > y_k \geq y_{k+1} = z_{k+1}; \end{aligned}$$

that is, the components of \mathbf{z} are arranged in order like those of \mathbf{x} and \mathbf{y} . Since \mathbf{y} majorizes \mathbf{x} , we find that

$$\sum_{i=1}^l z_i = \sum_{i=1}^l y_i \geq \sum_{i=1}^l x_i, \quad (10.3)$$

for $l = 1, \dots, j - 1$ and $l = k, \dots, m$, with equality when $l = m$. Using (10.3) when $l = j - 1$ and the fact that $z_j \geq x_j$, we have

$$\sum_{i=1}^l z_i \geq \sum_{i=1}^l x_i \quad (10.4)$$

when $l = j$, while (10.2) then guarantees that (10.4) holds for $l = j + 1, \dots, k - 1$, and so we have established that $\mathbf{x} \prec \mathbf{z}$. All that remains is to show that the number of nonzero differences $z_i - x_i$ is less than h . Since $z_j = x_j$ if $\delta = y_j - x_j$ and $z_k = x_k$ if $\delta = x_k - y_k$, it follows that the number of nonzero differences $z_i - x_i$ is $h - 1$ or $h - 2$. \square

The following is an immediate consequence of the proof of Theorem 10.2.

Corollary 10.2.1 If $\mathbf{x} \prec \mathbf{y}$, then there exist a finite number of T -transforms, T_1, \dots, T_n , such that $\mathbf{x} = P\mathbf{y}$, where $P = T_1 \cdots T_n$.

Our next result shows that there is a majorization relationship between the vector of diagonal elements of a symmetric matrix and the vector of its eigenvalues.

Theorem 10.3 Let A be an $m \times m$ symmetric matrix with diagonal elements a_{11}, \dots, a_{mm} and eigenvalues $\lambda_1(A) \geq \dots \geq \lambda_m(A)$. If $\mathbf{a} = (a_{11}, \dots, a_{mm})'$ and $\boldsymbol{\lambda} = (\lambda_1(A), \dots, \lambda_m(A))'$, then

$$\mathbf{a} \prec \boldsymbol{\lambda}.$$

Proof. We may assume that $a_{11} \geq \dots \geq a_{mm}$ since, if this is not the case, there is a permutation matrix P for which PAP' has nonincreasing diagonal elements and the same eigenvalues as A . For $k = 1, \dots, m-1$, let A_k be the leading $k \times k$ principal submatrix of A . From Theorem 3.20, $\lambda_i(A_k) \leq \lambda_i(A)$ for $i = 1, \dots, k$, so we have

$$\sum_{i=1}^k a_{ii} = \text{tr}(A_k) = \sum_{i=1}^k \lambda_i(A_k) \leq \sum_{i=1}^k \lambda_i(A).$$

The equality part of Definition 10.1 follows immediately from the fact that $\text{tr}(A) = \sum_{i=1}^m \lambda_i(A)$. \square

Our next result gives the converse of Theorem 10.3

Theorem 10.4 Suppose \mathbf{x} and \mathbf{y} are $m \times 1$ vectors and $\mathbf{x} \prec \mathbf{y}$. Then there exists an $m \times m$ symmetric matrix A with diagonal elements x_1, \dots, x_m and eigenvalues y_1, \dots, y_m .

Proof. Our proof follows that of Horn and Johnson (2013). Without loss of generality we assume that $x_1 \geq \dots \geq x_m$ and $y_1 \geq \dots \geq y_m$. Note that since $\mathbf{x} \prec \mathbf{y}$, it follows that $y_1 \geq x_1 \geq x_m \geq y_m$. Thus, the result immediately follows if $y_1 = y_m$, since in this case $\mathbf{x} = \mathbf{y} = y_1 \mathbf{1}_m$ and the required matrix is given by $A = y_1 I_m$. For the remainder of the proof we assume that $y_1 > y_m$. We first consider the case in which $m = 2$, so that $y_1 \geq x_1 \geq x_2 \geq y_2$ with $y_1 > y_2$. Now the matrix

$$P = \frac{1}{\sqrt{y_1 - y_2}} \begin{bmatrix} \sqrt{y_1 - x_2} & \sqrt{x_2 - y_2} \\ -\sqrt{x_2 - y_2} & \sqrt{y_1 - x_2} \end{bmatrix}$$

is orthogonal, and so if $\Lambda = \text{diag}(y_1, y_2)$, then $A = P\Lambda P'$ is a symmetric matrix with eigenvalues y_1 and y_2 . The diagonal elements of A reduce to $y_1 + y_2 - x_2$ and x_2 , but since we must have $x_1 + x_2 = y_1 + y_2$, we see the diagonal elements are in fact x_1 and x_2 as required. This proves the result for $m = 2$, and we will prove the result

for larger m by induction. Suppose that $m > 2$ and assume that the result holds for values less than or equal to $m - 1$. Let k be the largest index for which $y_k \geq x_1$. Such a k must exist since $y_1 \geq x_1$. Note that if $k = m$, then $y_m \geq x_1 \geq x_m \geq y_m$ and so $x_1 = \cdots = x_m = y_m$. But then $y_i - x_1 \geq y_m - x_1 \geq 0$ and

$$\sum_{i=1}^m (y_i - x_1) = \sum_{i=1}^m (y_i - x_i) = \sum_{i=1}^m y_i - \sum_{i=1}^m x_i = 0,$$

imply each $y_i = x_1$, and this contradicts the assumption that $y_1 > y_m$. Thus, $k \leq m - 1$ and so $y_k \geq x_1 > y_{k+1} \geq y_m$. This and the identity $x_1 + (y_k + y_{k+1} - x_1) = y_k + y_{k+1}$ show that $\mathbf{u} = (x_1, y_k + y_{k+1} - x_1)' \prec \mathbf{v} = (y_k, y_{k+1})'$. Since $y_k > y_{k+1}$, we can use the construction for the case $m = 2$ to obtain an orthogonal matrix P_1 such that $P_1 D_{\mathbf{v}} P_1'$ has diagonal elements $u_1 = x_1$ and $u_2 = y_k + y_{k+1} - x_1$. Letting $\mathbf{z}' = (u_2, \mathbf{z}_2)'$, where \mathbf{z}_2 is the $(m - 2) \times 1$ vector obtained from \mathbf{y} by removing y_k and y_{k+1} , we have

$$\begin{bmatrix} P_1 & (0) \\ (0) & I_{m-2} \end{bmatrix} \begin{bmatrix} D_{\mathbf{v}} & (0) \\ (0) & D_{\mathbf{z}_2} \end{bmatrix} \begin{bmatrix} P_1' & (0) \\ (0) & I_{m-2} \end{bmatrix} = \begin{bmatrix} x_1 & \mathbf{a}' \\ \mathbf{a} & D_{\mathbf{z}} \end{bmatrix}$$

for some $(m - 1) \times 1$ vector \mathbf{a} . The result will then follow if we can find an $(m - 1) \times (m - 1)$ orthogonal matrix P_2 such that $P_2 D_{\mathbf{z}} P_2'$ has diagonal elements x_2, \dots, x_m . Since we are assuming the result holds for $m - 1$, the proof will be complete if we can show $\mathbf{w} = (x_2, \dots, x_m)' \prec \mathbf{z}$. Note that

$$y_k = u_2 + (x_1 - y_{k+1}) > u_2 = y_{k+1} + (y_k - x_1) \geq y_{k+1}. \quad (10.5)$$

First suppose that $k = 1$. It follows from (10.5) that the components of \mathbf{z} , like those of \mathbf{w} , are nondecreasing. Then for $h = 1, \dots, m - 1$,

$$\begin{aligned} \sum_{i=1}^h z_i &= u_2 + \sum_{i=2}^h y_{i+1} = y_1 + y_2 - x_1 + \sum_{i=3}^{h+1} y_i \\ &= \sum_{i=1}^{h+1} y_i - x_1 \geq \sum_{i=1}^{h+1} x_i - x_1 = \sum_{i=2}^{h+1} x_i = \sum_{i=1}^h w_i, \end{aligned}$$

with equality if $h = m - 1$. Now suppose $k > 1$, in which case the ordered components of \mathbf{z} are $y_1 \geq \cdots \geq y_{k-1} \geq u_2 \geq y_{k+2} \geq \cdots \geq y_m$. Then

$$\sum_{i=1}^h z_{[i]} = \sum_{i=1}^h y_i \geq \sum_{i=1}^h x_i \geq \sum_{i=1}^h x_{i+1} = \sum_{i=1}^h w_i,$$

for $h = 1, \dots, k-1$,

$$\begin{aligned} \sum_{i=1}^k z_{[i]} &= \sum_{i=1}^{k-1} y_i + u_2 = \sum_{i=1}^{k-1} y_i + y_k + y_{k+1} - x_1 = \sum_{i=1}^{k+1} y_i - x_1 \\ &\geq \sum_{i=1}^{k+1} x_i - x_1 = \sum_{i=2}^{k+1} x_i = \sum_{i=1}^k w_i, \end{aligned} \quad (10.6)$$

and

$$\begin{aligned} \sum_{i=1}^h z_{[i]} &= \sum_{i=1}^{k-1} y_i + u_2 + \sum_{i=k+1}^h y_{i+1} = \sum_{i=1}^{k-1} y_i + y_k + y_{k+1} - x_1 + \sum_{i=k+2}^{h+1} y_i \\ &= \sum_{i=1}^{h+1} y_i - x_1 \geq \sum_{i=1}^{h+1} x_i - x_1 = \sum_{i=2}^{h+1} x_i = \sum_{i=1}^h w_i, \end{aligned} \quad (10.7)$$

for $h = k+1, \dots, m-1$. Note that we have equality in (10.6) when $k = m-1$ and equality in (10.7) when $h = m-1$, so we have shown that $\mathbf{w} \prec \mathbf{z}$. \square

Our next result gives a majorization relationship between the sum of the two vectors of eigenvalues of two symmetric matrices and the vector of eigenvalues of the sum of the two matrices.

Theorem 10.5 Let A and B be $m \times m$ symmetric matrices. If the i th components of \mathbf{a} , \mathbf{b} and \mathbf{c} are $\lambda_i(A)$, $\lambda_i(B)$ and $\lambda_i(A+B)$, respectively, then $\mathbf{c} \prec (\mathbf{a} + \mathbf{b})$.

Proof. The result can be proven using the extremal properties of eigenvalues. If P is an $m \times k$ semiorthogonal matrix, then for $k = 1, \dots, m$, we have

$$\begin{aligned} \sum_{i=1}^k c_i &= \sum_{i=1}^k \lambda_i(A+B) = \max_{P'P=I_k} \operatorname{tr}\{P'(A+B)P\} \\ &= \max_{P'P=I_k} \{\operatorname{tr}(P'AP) + \operatorname{tr}(P'BP)\} \\ &\leq \max_{P'P=I_k} \operatorname{tr}(P'AP) + \max_{P'P=I_k} \operatorname{tr}(P'BP) \\ &= \sum_{i=1}^k \lambda_i(A) + \sum_{i=1}^k \lambda_i(B) = \sum_{i=1}^k (a_i + b_i). \end{aligned}$$

We have equality when $k = m$ since $\operatorname{tr}(A+B) = \operatorname{tr}(A) + \operatorname{tr}(B)$. \square

The result given in Theorem 10.5 holds when the components of \mathbf{a} and \mathbf{b} have been ordered in the same manner. When one has nondecreasing components and the other has nonincreasing components, we get the opposite majorization relationship.

Theorem 10.6 Let A and B be $m \times m$ symmetric matrices. If the i th components of \mathbf{a} and \mathbf{c} are $\lambda_i(A)$ and $\lambda_i(A+B)$, respectively, and the i th component of \mathbf{b} is $\lambda_{m-i+1}(B)$, then $(\mathbf{a} + \mathbf{b}) \prec \mathbf{c}$.

Proof. If P is an $m \times k$ semiorthogonal matrix, then for $k = 1, \dots, m$, we have

$$\begin{aligned} \sum_{i=1}^k c_i &= \sum_{i=1}^k \lambda_i(A+B) = \max_{P'P=I_k} \operatorname{tr}\{P'(A+B)P\} \\ &= \max_{P'P=I_k} \{\operatorname{tr}(P'AP) + \operatorname{tr}(P'BP)\} \\ &\geq \max_{P'P=I_k} \{\operatorname{tr}(P'AP) + \min_{P'P=I_k} \operatorname{tr}(P'BP)\} \\ &= \max_{P'P=I_k} \left\{ \operatorname{tr}(P'AP) + \sum_{i=1}^k \lambda_{m-i+1}(B) \right\} \\ &= \sum_{i=1}^k \lambda_i(A) + \sum_{i=1}^k \lambda_{m-i+1}(B) = \sum_{i=1}^k (a_i + b_i). \end{aligned}$$

We have equality when $k = m$ since $\operatorname{tr}(A+B) = \operatorname{tr}(A) + \operatorname{tr}(B)$. □

Our next result gives upper and lower bounds for $\operatorname{tr}(AB)$ in terms of the eigenvalues of A and B .

Theorem 10.7 If A and B are $m \times m$ symmetric matrices, then

$$\sum_{i=1}^m \lambda_i(A) \lambda_{m-i+1}(B) \leq \operatorname{tr}(AB) \leq \sum_{i=1}^m \lambda_i(A) \lambda_i(B). \quad (10.8)$$

Proof. Let $A = P\Lambda P'$ be the spectral decomposition of A , so that P is an $m \times m$ orthogonal matrix and $\Lambda = \operatorname{diag}(\lambda_1(A), \dots, \lambda_m(A))$. Then if $C = P'BP$,

$$\begin{aligned} \operatorname{tr}(AB) &= \operatorname{tr}(P\Lambda P'B) = \operatorname{tr}(\Lambda P'BP) \\ &= \operatorname{tr}(\Lambda C) = \sum_{i=1}^m \lambda_i(A) c_{ii}. \end{aligned} \quad (10.9)$$

Since the eigenvalues of C are the same as those of B , from Theorem 10.3, $\mathbf{c} \prec \mathbf{b}$, where $\mathbf{c} = (c_{11}, \dots, c_{mm})'$ and $\mathbf{b} = (\lambda_1(B), \dots, \lambda_m(B))'$. Using this along with the fact that $\lambda_i(A) - \lambda_{i+1}(A)$ is nonnegative, we have

$$\begin{aligned}
 \sum_{i=1}^m \lambda_i(A) c_{ii} &= \lambda_1(A) c_{11} + \sum_{i=2}^m \lambda_i(A) \left(\sum_{j=1}^i c_{jj} - \sum_{j=1}^{i-1} c_{jj} \right) \\
 &= \lambda_1(A) c_{11} + \sum_{i=2}^m \lambda_i(A) \sum_{j=1}^i c_{jj} - \sum_{i=1}^{m-1} \lambda_{i+1}(A) \sum_{j=1}^i c_{jj} \\
 &= \sum_{i=1}^{m-1} \{ \lambda_i(A) - \lambda_{i+1}(A) \} \sum_{j=1}^i c_{jj} + \lambda_m(A) \sum_{j=1}^m c_{jj} \\
 &\leq \sum_{i=1}^{m-1} \{ \lambda_i(A) - \lambda_{i+1}(A) \} \sum_{j=1}^i \lambda_j(B) + \lambda_m(A) \sum_{j=1}^m c_{jj} \\
 &= \sum_{i=1}^{m-1} \{ \lambda_i(A) - \lambda_{i+1}(A) \} \sum_{j=1}^i \lambda_j(B) + \lambda_m(A) \sum_{j=1}^m \lambda_j(B) \\
 &= \lambda_1(A) \lambda_1(B) + \sum_{i=2}^m \lambda_i(A) \sum_{j=1}^i \lambda_j(B) \\
 &\quad - \sum_{i=1}^{m-1} \lambda_{i+1}(A) \sum_{j=1}^i \lambda_j(B) \\
 &= \lambda_1(A) \lambda_1(B) + \sum_{i=2}^m \lambda_i(A) \left(\sum_{j=1}^i \lambda_j(B) - \sum_{j=1}^{i-1} \lambda_j(B) \right) \\
 &= \sum_{i=1}^m \lambda_i(A) \lambda_i(B). \tag{10.10}
 \end{aligned}$$

Combining (10.9) and (10.10) yields the upper bound in (10.8). Applying this upper bound to A and $-B$, we get

$$-\operatorname{tr}(AB) \leq \sum_{i=1}^m \lambda_i(A) \lambda_i(-B) = - \sum_{i=1}^m \lambda_i(A) \lambda_{m-i+1}(B),$$

which is equivalent to the lower bound in (10.8). \square

Theorem 10.7 can be used to obtain a more general result regarding a partial sum of eigenvalues of a matrix product.

Theorem 10.8 Let A and B be $m \times m$ symmetric matrices with A being nonnegative definite. Then

$$\sum_{i=1}^k \lambda_i(AB) \geq \sum_{i=1}^k \lambda_i(A) \lambda_{m-i+1}(B),$$

for $k = 1, \dots, m$.

Proof. Let $A = PDP'$ be the spectral decomposition of A , where the diagonal matrix is given by $D = \text{diag}(\lambda_1(A), \dots, \lambda_m(A))$, and let D_* be the $m \times m$ diagonal matrix $D_* = \text{diag}(\lambda_1(A), \dots, \lambda_k(A), 0, \dots, 0)$. It follows from Problem 3.12 that the eigenvalues of AB are the same as those of $D^{1/2}P'BP D^{1/2}$. Using this and Problem 3.57, we have

$$\begin{aligned} \sum_{i=1}^k \lambda_i(AB) &= \sum_{i=1}^k \lambda_i(D^{1/2}P'BP D^{1/2}) \\ &= \max_{C'C=I_k} \text{tr}(C'D^{1/2}P'BP D^{1/2}C) \\ &\geq \text{tr} \left\{ [I_k, (0)] D^{1/2}P'BP D^{1/2} \begin{bmatrix} I_k \\ (0) \end{bmatrix} \right\} \\ &= \text{tr} \left(D^{1/2} \begin{bmatrix} I_k & (0) \\ (0) & (0) \end{bmatrix} D^{1/2}P'BP \right) = \text{tr}(D_*P'BP) \\ &\geq \sum_{i=1}^m \lambda_i(D_*) \lambda_{m-i+1}(P'BP) = \sum_{i=1}^k \lambda_i(A) \lambda_{m-i+1}(B), \end{aligned}$$

where the last inequality follows from Theorem 10.7. \square

Majorization can be combined with order-preserving functions to develop numerous useful inequalities. A real-valued function $\phi(\mathbf{x})$ defined for all $\mathbf{x} \in S$, where S is a subset of R^m , is said to be order-preserving if $\mathbf{x} \prec \mathbf{y}$ implies that $\phi(\mathbf{x}) \leq \phi(\mathbf{y})$. We will focus on functions $\phi(\mathbf{x})$ that have the form $\phi(\mathbf{x}) = \sum_{i=1}^m g(x_i)$, where $g(x)$ is a real-valued scalar function. Our final result of this section gives a sufficient condition for this type of function to be order-preserving.

Theorem 10.9 Let g be a real-valued convex function defined on an interval S in R . If $\mathbf{x} \prec \mathbf{y}$, then

$$\sum_{i=1}^m g(x_i) \leq \sum_{i=1}^m g(y_i). \quad (10.11)$$

In addition, if g is strictly convex on S , then we have equality in (10.11) if and only if $x_{[i]} = y_{[i]}$ for $i = 1, \dots, m$.

Proof. We assume without loss of generality that $x_1 \geq \cdots \geq x_m$ and $y_1 \geq \cdots \geq y_m$. It follows from Corollary 10.2.1 that there are $m \times 1$ vectors x_1, \dots, x_{n-1} such that

$$\mathbf{x} = \mathbf{x}_0 \prec \mathbf{x}_1 \prec \cdots \prec \mathbf{x}_{n-1} \prec \mathbf{x}_n = \mathbf{y}$$

and x_{i-1} and x_i have at most 2 components that differ. Consequently, it will suffice to prove the result for $m = 2$. Since $\mathbf{x} \prec \mathbf{y}$, there exists a doubly stochastic matrix P such that $\mathbf{x} = P\mathbf{y}$. Every 2×2 doubly stochastic matrix has the form

$$P = \begin{bmatrix} \alpha & 1 - \alpha \\ 1 - \alpha & \alpha \end{bmatrix}$$

for some $\alpha \in [0, 1]$, so we have $x_1 = \alpha y_1 + (1 - \alpha)y_2$ and $x_2 = (1 - \alpha)y_1 + \alpha y_2$. Then since g is convex

$$\begin{aligned} g(x_1) + g(x_2) &= g(\alpha y_1 + (1 - \alpha)y_2) + g((1 - \alpha)y_1 + \alpha y_2) \\ &\leq \{\alpha g(y_1) + (1 - \alpha)g(y_2)\} + \{(1 - \alpha)g(y_1) + \alpha g(y_2)\} \\ &= g(y_1) + g(y_2), \end{aligned}$$

which establishes (10.11). Clearly, if $\mathbf{x} = \mathbf{y}$, we have equality in (10.11). Next assume we have equality in (10.11) and note that this requires

$$g(\alpha y_1 + (1 - \alpha)y_2) = \alpha g(y_1) + (1 - \alpha)g(y_2),$$

and

$$g((1 - \alpha)y_1 + \alpha y_2) = (1 - \alpha)g(y_1) + \alpha g(y_2).$$

But if g is strictly convex these identities only hold if $x_1 = y_1$ and $x_2 = y_2$. □

Example 10.1 As a simple example of Theorem 10.9, we consider the strictly convex function $g(x) = (x - a)^2$, where a is a real number. Then it follows that if $\mathbf{x} \prec \mathbf{y}$,

$$\sum_{i=1}^m (x_i - a)^2 \leq \sum_{i=1}^m (y_i - a)^2,$$

with equality if and only if $x_{[i]} = y_{[i]}$, for $i = 1, \dots, m$. In particular, if we choose $a = \bar{x} = m^{-1} \sum_{i=1}^m x_i$, so that $a = \bar{y}$ also, we have

$$s_x^2 = \frac{1}{m-1} \sum_{i=1}^m (x_i - \bar{x})^2 \leq \frac{1}{m-1} \sum_{i=1}^m (y_i - \bar{y})^2 = s_y^2 \quad (10.12)$$

if $\mathbf{x} \prec \mathbf{y}$. Since $\mathbf{x} \prec \mathbf{y}$ means that the components of \mathbf{x} are less spread out than those of \mathbf{y} in some sense, namely as indicated by Definition 10.1, the inequality in (10.12) is not surprising.

10.3 CAUCHY-SCHWARZ INEQUALITIES

We will start by restating the standard Cauchy-Schwarz inequality that was presented in Section 2.2 and give an alternative proof.

Theorem 10.10 Suppose that \mathbf{x} and \mathbf{y} are $m \times 1$ vectors. Then

$$(\mathbf{x}'\mathbf{y})^2 \leq (\mathbf{x}'\mathbf{x})(\mathbf{y}'\mathbf{y}), \quad (10.13)$$

with equality if and only if one of the vectors is a scalar multiple of the other.

Proof. If either vector is the null vector, then (10.13) holds with equality and that null vector equals zero times the other vector. Otherwise, note that

$$(\mathbf{x}'\mathbf{x})(\mathbf{y}'\mathbf{y}) - (\mathbf{x}'\mathbf{y})^2 = (\mathbf{x}'\mathbf{x})\mathbf{y}'A\mathbf{y}, \quad (10.14)$$

where $A = I_m - (\mathbf{x}'\mathbf{x})^{-1}\mathbf{x}\mathbf{x}'$. Clearly, A has eigenvalues of 1 and 0 with multiplicities of $m - 1$ and 1. This guarantees that (10.14) is nonnegative and so (10.13) holds. Eigenvectors of A corresponding to the 0 eigenvalue are of the form $c\mathbf{x}$ for $c \neq 0$, and so (10.14) equals 0, and hence (10.13) holds with equality, if and only if \mathbf{y} has this same form. \square

A simple extension of the standard Cauchy-Schwarz inequality is given next.

Theorem 10.11 If \mathbf{x} and \mathbf{y} are $m \times 1$ vectors and A is an $m \times m$ positive definite matrix, then

$$(\mathbf{x}'\mathbf{y})^2 \leq (\mathbf{x}'A\mathbf{x})(\mathbf{y}'A^{-1}\mathbf{y}),$$

with equality if and only if one of the vectors, $A\mathbf{x}$ and \mathbf{y} , is a scalar multiple of the other.

Proof. Since A is positive definite, there exists a nonsingular matrix T such that $A = T'T$ and $A^{-1} = T^{-1}T^{-1'}$. Defining $\mathbf{u} = T\mathbf{x}$ and $\mathbf{v} = T^{-1'}\mathbf{y}$, we find that

$$\begin{aligned} (\mathbf{x}'\mathbf{y})^2 &= (\mathbf{x}'T'T^{-1'}\mathbf{y})^2 = (\mathbf{u}'\mathbf{v})^2 \leq (\mathbf{u}'\mathbf{u})(\mathbf{v}'\mathbf{v}) \\ &= (\mathbf{x}'T'T\mathbf{x})(\mathbf{y}'T^{-1}T^{-1'}\mathbf{y}) = (\mathbf{x}'A\mathbf{x})(\mathbf{y}'A^{-1}\mathbf{y}), \end{aligned}$$

where the inequality follows from Theorem 10.10. We have equality if and only if one of the vectors, $\mathbf{u} = T\mathbf{x}$ and $\mathbf{v} = T^{-1'}\mathbf{y}$, is a scalar multiple of the other. Premultiplying by T' , we see this is equivalent to saying that one of the vectors, $A\mathbf{x}$ and \mathbf{y} , is a scalar multiple of the other. \square

The rest of this section is devoted to matrix versions of the Cauchy-Schwarz inequality. We first consider one involving the trace.

Theorem 10.12 If A and B are both $m \times n$ matrices, then

$$\{\operatorname{tr}(A'B)\}^2 \leq \{\operatorname{tr}(A'A)\}\{\operatorname{tr}(B'B)\},$$

with equality if and only if one of the matrices is a scalar multiple of the other.

Proof. Using Theorem 8.10, we can write $\operatorname{tr}(A'B) = \mathbf{x}'\mathbf{y}$, $\operatorname{tr}(A'A) = \mathbf{x}'\mathbf{x}$ and $\operatorname{tr}(B'B) = \mathbf{y}'\mathbf{y}$, where $\mathbf{x} = \operatorname{vec}(A)$ and $\mathbf{y} = \operatorname{vec}(B)$. The result then follows immediately from Theorem 10.10. \square

Before giving a general Cauchy-Schwarz inequality involving determinants, we first consider the following special case.

Theorem 10.13 Suppose that P and Q are both $m \times n$ semiorthogonal matrices with $n \leq m$. Then

$$|P'Q|^2 \leq 1,$$

with equality if and only if $PP' = QQ'$.

Proof. We have

$$|P'Q|^2 = |Q'PP'Q| \leq |Q'PP'Q + Q'(I_m - PP')Q| = |I_n| = 1, \quad (10.15)$$

where the inequality follows from Theorem 4.17 since $Q'PP'Q$ and $Q'(I_m - PP')Q$ are both nonnegative definite. From that same theorem, we see that we have equality in (10.15) if and only if $Q'(I_m - PP')Q = (0)$ which is equivalent to the condition $PP' = QQ'$. \square

Next we have the generalization of Theorem 10.13 to arbitrary real matrices.

Theorem 10.14 Suppose that both A and B are $m \times n$ matrices. Then

$$|A'B|^2 \leq |A'A| |B'B|, \quad (10.16)$$

with equality if and only if $\operatorname{rank}(A) < n$ or $\operatorname{rank}(B) < n$, or $B = AC$ for some nonsingular matrix C .

Proof. Clearly the inequality holds when $|A'B| = 0$ and, in this case, equality holds if and only if $\operatorname{rank}(A) < n$ or $\operatorname{rank}(B) < n$. For the remainder of the proof we

assume $|A'B| \neq 0$. Using the singular value decomposition we can write A and B as $A = P_1 D_1 Q_1$ and $B = P_2 D_2 Q_2$, where the $m \times n$ matrix P_i and $n \times n$ matrix Q_i satisfy $P_i' P_i = Q_i' Q_i = I_n$, and D_i is an $n \times n$ diagonal matrix with positive diagonal elements. It then follows that

$$|A'B|^2 = |Q_1' D_1 P_1' P_2 D_2 Q_2|^2 = |D_1|^2 |D_2|^2 |P_1' P_2|^2,$$

while $|A'A| = |D_1|^2$ and $|B'B| = |D_2|^2$. Thus, (10.16) follows directly from Theorem 10.13. Also from Theorem 10.13, we have equality if and only if $P_1 P_1' = P_2 P_2'$, and since this is equivalent to A and B having the same column space, the proof is complete. \square

10.4 HÖLDER'S INEQUALITY

We start with one version of an inequality known as Hölder's inequality. If \mathbf{x} and \mathbf{y} are $m \times 1$ vectors with nonnegative components and α is a scalar satisfying $0 < \alpha < 1$, then

$$\sum_{i=1}^m x_i^\alpha y_i^{1-\alpha} \leq \left(\sum_{i=1}^m x_i \right)^\alpha \left(\sum_{i=1}^m y_i \right)^{1-\alpha}, \quad (10.17)$$

with equality if and only if one of the vectors is a scalar multiple of the other. A proof of this result, along with some extensions to more than two vectors, can be found in Hardy et al. (1952). We will first use Hölder's inequality to prove the following result.

Theorem 10.15 Suppose \mathbf{a} is an $m \times 1$ vector with nonnegative components and $\frac{1}{p} + \frac{1}{q} = 1$, where $p > 1$. Then

$$\sum_{i=1}^m a_i b_i \leq \left(\sum_{i=1}^m a_i^p \right)^{1/p} \left(\sum_{i=1}^m b_i^q \right)^{1/q}, \quad (10.18)$$

for every $m \times 1$ vector \mathbf{b} having nonnegative components and satisfying $\sum_{i=1}^m b_i^q = 1$. We have equality in (10.18) if and only if $\mathbf{a} = \mathbf{0}$ or

$$b_i^q = \frac{a_i^p}{\sum_{j=1}^m a_j^p},$$

for $i = 1, \dots, m$.

Proof. Clearly both sides of (10.18) reduce to 0 when $\mathbf{a} = \mathbf{0}$, so for the remainder of the proof we assume $\mathbf{a} \neq \mathbf{0}$. Let $x_i = a_i^p$, $y_i = b_i^q$, and $\alpha = 1/p$ so that $a_i = x_i^\alpha$ and $b_i = y_i^{1-\alpha}$. Then using (10.17), we have

$$\begin{aligned}
\sum_{i=1}^m a_i b_i &= \sum_{i=1}^m x_i^\alpha y_i^{1-\alpha} \leq \left(\sum_{i=1}^m x_i \right)^\alpha \left(\sum_{i=1}^m y_i \right)^{1-\alpha} \\
&= \left(\sum_{i=1}^m a_i^p \right)^{1/p} \left(\sum_{i=1}^m b_i^q \right)^{1/q} = \left(\sum_{i=1}^m a_i^p \right)^{1/p}.
\end{aligned}$$

Further, we have equality if and only if $\mathbf{x} = (a_1^p, \dots, a_m^p)'$ and $\mathbf{y} = (b_1^q, \dots, b_m^q)'$ are scalar multiples of one another which, due to the constraint $\sum_{i=1}^m b_i^q = 1$, implies that

$$b_i^q = \frac{a_i^p}{\sum_{j=1}^m a_j^p}.$$

□

In this section, we will look at some matrix versions of Hölder's inequality. Our first of these results involves determinants and is due to Fan (1950).

Theorem 10.16 Suppose that A and B are $m \times m$ nonnegative definite matrices and α is a scalar satisfying $0 < \alpha < 1$. Then

$$|A|^\alpha |B|^{1-\alpha} \leq |\alpha A + (1-\alpha)B|, \quad (10.19)$$

with equality if and only if $A = B$ or $\alpha A + (1-\alpha)B$ is singular.

Proof. Since $\alpha A + (1-\alpha)B$ is also nonnegative definite, (10.19) clearly holds when A or B is singular, with equality if and only if $\alpha A + (1-\alpha)B$ is also singular. For the remainder of the proof we assume that both A and B are positive definite. Using Theorem 4.14, we can write $A = T\Lambda T'$ and $B = TT'$, where T is a nonsingular matrix, $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$, and $\lambda_1, \dots, \lambda_m$ are the eigenvalues of $B^{-1}A$. Thus, the proof will be complete if we can show that

$$|\Lambda|^\alpha = \prod_{i=1}^m \lambda_i^\alpha \leq |\alpha \Lambda + (1-\alpha)I_m| = \prod_{i=1}^m (\alpha \lambda_i + 1 - \alpha),$$

with equality if and only if $\Lambda = I_m$. This result is easily confirmed by showing the function $g(\lambda) = \alpha \lambda + 1 - \alpha - \lambda^\alpha$ is minimized at $\lambda = 1$ when $0 < \alpha < 1$. □

For nonnegative definite matrices A and B , we will next give an upper bound for $\text{tr}(A^\alpha B^{1-\alpha})$, where $0 < \alpha < 1$. Here the α th power of the nonnegative definite matrix A is defined to be $A^\alpha = X\Lambda^\alpha X'$, where $A = X\Lambda X'$ represents the spectral decomposition of A and $\Lambda^\alpha = \text{diag}(\lambda_1^\alpha, \dots, \lambda_m^\alpha)$. Before establishing

the bound, we will need a preliminary result. Both of these results are due to Magnus (1987).

Theorem 10.17 Suppose A is an $m \times m$ nonnull nonnegative definite matrix and $\frac{1}{p} + \frac{1}{q} = 1$, where $p > 1$. Then

$$\operatorname{tr}(AB) \leq \{\operatorname{tr}(A^p)\}^{1/p} \quad (10.20)$$

for every $m \times m$ nonnegative definite matrix B satisfying $\operatorname{tr}(B^q) = 1$. We have equality in (10.20) if and only if

$$B^q = \{\operatorname{tr}(A^p)\}^{-1} A^p.$$

Proof. Let $B = P\Lambda P'$ be the spectral decomposition of B , and note that from the conditions of the theorem, $\sum_{i=1}^m \lambda_i^q = 1$, where $\lambda_1 \geq \dots \geq \lambda_m$ are the eigenvalues of B . Put $C = P'AP$ so that

$$\operatorname{tr}(AB) = \operatorname{tr}(AP\Lambda P') = \operatorname{tr}(C\Lambda) = \sum_{i=1}^m c_{ii}\lambda_i,$$

where c_{ii} is the i th diagonal element of C . Applying Theorem 10.15, we have

$$\operatorname{tr}(AB) \leq \left(\sum_{i=1}^m c_{ii}^p \right)^{1/p}. \quad (10.21)$$

Since $(c_{11}, \dots, c_{mm})' \prec (\gamma_1, \dots, \gamma_m)'$, where $\gamma_1 \geq \dots \geq \gamma_m$ are the eigenvalues of A and C , and $g(x) = x^p$ is a strictly convex function on $[0, \infty)$ when $p > 1$, an application of Theorem 10.9 leads to

$$\sum_{i=1}^m c_{ii}^p \leq \sum_{i=1}^m \gamma_i^p = \operatorname{tr}(C^p) = \operatorname{tr}(A^p). \quad (10.22)$$

Combining (10.21) and (10.22) immediately yields (10.20). From Theorem 10.15, equality occurs in (10.21) if and only if

$$\lambda_i^q = \frac{c_{ii}^p}{\sum_{j=1}^m c_{jj}^p},$$

for $i = 1, \dots, m$, and by Theorem 10.9, we have equality in (10.22) if and only if c_{11}, \dots, c_{mm} are the eigenvalues of C , that is, $C = \Gamma = \operatorname{diag}(\gamma_1, \dots, \gamma_m)$ and $A = P\Gamma P'$. Thus, we have equality in (10.20) if and only if

$$\begin{aligned}
B^q &= P \Lambda^q P' = \left(\sum_{j=1}^m c_{jj}^p \right)^{-1} P \operatorname{diag}(c_{11}^p, \dots, c_{mm}^p) P' \\
&= \left(\sum_{j=1}^m \gamma_j^p \right)^{-1} P \Gamma^p P' = \{\operatorname{tr}(A^p)\}^{-1} A^p,
\end{aligned}$$

and so the proof is complete. \square

Now we can use Theorem 10.17 to establish a matrix version of Hölder's inequality involving traces.

Theorem 10.18 Suppose A and B are $m \times m$ nonnull nonnegative definite matrices and α is a scalar satisfying $0 < \alpha < 1$. Then

$$\operatorname{tr}(A^\alpha B^{1-\alpha}) \leq \{\operatorname{tr}(A)\}^\alpha \{\operatorname{tr}(B)\}^{1-\alpha}, \quad (10.23)$$

with equality if and only if $B = cA$ for some positive scalar c .

Proof. Let $p = \alpha^{-1}$ and $q = (1 - \alpha)^{-1}$, so that $\frac{1}{p} + \frac{1}{q} = 1$ and $p > 1$. Define

$$C = \frac{B^{1/q}}{\{\operatorname{tr}(B)\}^{1/q}}$$

so that $\operatorname{tr}(C^q) = 1$, and then by Theorem 10.17,

$$\operatorname{tr}(A^{1/p} C) \leq [\operatorname{tr}\{(A^{1/p})^p\}]^{1/p} = \{\operatorname{tr}(A)\}^{1/p} = \{\operatorname{tr}(A)\}^\alpha.$$

This yields (10.23) since

$$\operatorname{tr}(A^{1/p} C) = \frac{\operatorname{tr}(A^{1/p} B^{1/q})}{\{\operatorname{tr}(B)\}^{1/q}} = \frac{\operatorname{tr}(A^\alpha B^{1-\alpha})}{\{\operatorname{tr}(B)\}^{1-\alpha}}.$$

According to Theorem 10.17, we have equality in (10.23) if and only if

$$C^q = \frac{(A^{1/p})^p}{\operatorname{tr}\{(A^{1/p})^p\}} = \frac{A}{\operatorname{tr}(A)},$$

or equivalently $B = \{\operatorname{tr}(B)/\operatorname{tr}(A)\}A$. \square

10.5 MINKOWSKI'S INEQUALITY

Minkowski's inequality is another well-known classical inequality. If \mathbf{x} and \mathbf{y} are $m \times 1$ vectors with nonnegative components and $p > 1$, then

$$\left\{ \sum_{i=1}^m (x_i + y_i)^p \right\}^{1/p} \leq \left(\sum_{i=1}^m x_i^p \right)^{1/p} + \left(\sum_{i=1}^m y_i^p \right)^{1/p}, \quad (10.24)$$

with equality if and only if \mathbf{x} and \mathbf{y} are linearly dependent. A proof of this result, along with some extensions can be found in Hardy et al. (1952). A Minkowski inequality for products is given by

$$\left\{ \prod_{i=1}^m (x_i + y_i) \right\}^{1/m} \geq \left(\prod_{i=1}^m x_i \right)^{1/m} + \left(\prod_{i=1}^m y_i \right)^{1/m}, \quad (10.25)$$

where again we have equality if and only if \mathbf{x} and \mathbf{y} are linearly dependent. The inequality (10.25) is actually a special case of a generalized Hölder inequality; see, for instance, Magnus and Neudecker (1999).

In this section, we will obtain some matrix analogues of Minkowski's inequalities. Our first result, which involves determinants, can be viewed as a generalization of (10.25).

Theorem 10.19 Let A and B be $m \times m$ nonnull nonnegative definite matrices. Then

$$|A + B|^{1/m} \geq |A|^{1/m} + |B|^{1/m}, \quad (10.26)$$

with equality if and only if $A + B$ is singular or $A = cB$ for some $c > 0$.

Proof. Since $A + B$ is nonnegative definite, the inequality clearly holds when both A and B are singular with equality if and only if $A + B$ is also singular. For the remainder of the proof, we assume without loss of generality that B is positive definite. Applying Theorem 4.14, there exists a nonsingular matrix T such that $A = T\Lambda T'$ and $B = TT'$, where $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ has the eigenvalues of $T^{-1/2}AT^{-1/2}$ as its diagonal elements. Since $T^{-1/2}AT^{-1/2}$ is nonnegative definite,

$\lambda_i \geq 0$ for $i = 1, \dots, m$, and so using (10.25)

$$\begin{aligned}
 |A + B|^{1/m} &= |T\Lambda T' + TT'|^{1/m} = |T|^{2/m} |\Lambda + I_m|^{1/m} \\
 &= |T|^{2/m} \left\{ \prod_{i=1}^m (\lambda_i + 1) \right\}^{1/m} \\
 &\geq |T|^{2/m} \left\{ \left(\prod_{i=1}^m \lambda_i \right)^{1/m} + \left(\prod_{i=1}^m 1 \right)^{1/m} \right\} \\
 &= |T|^{2/m} (|\Lambda|^{1/m} + |I_m|^{1/m}) = |T\Lambda T'|^{1/m} + |TT'|^{1/m} \\
 &= |A|^{1/m} + |B|^{1/m},
 \end{aligned}$$

which establishes (10.26). We have equality if and only if $(\lambda_1, \dots, \lambda_m)'$ and $\mathbf{1}_m$ are linearly dependent; that is, $\lambda_1 = \dots = \lambda_m = c$, so that $\Lambda = cI_m$ and $A = T\Lambda T' = cTT' = cB$. \square

Our next result is a matrix version of (10.24) involving traces. The proof of this result, which utilizes Theorem 10.17, is due to Magnus (1987).

Theorem 10.20 Let A and B be $m \times m$ nonnull nonnegative definite matrices and suppose $p > 1$. Then

$$[\operatorname{tr}\{(A + B)^p\}]^{1/p} \leq \{\operatorname{tr}(A^p)\}^{1/p} + \{\operatorname{tr}(B^p)\}^{1/p},$$

with equality if and only if $B = cA$ for some $c > 0$.

Proof. It follows from Theorem 10.17 that for any $m \times m$ nonnull nonnegative definite matrix A

$$\max_{C \in \mathcal{S}} \operatorname{tr}(AC) = \{\operatorname{tr}(A^p)\}^{1/p},$$

where

$$\mathcal{S} = \{C : C \text{ is } m \times m \text{ nonnegative definite, } \operatorname{tr}(C^q) = 1\},$$

and q is such that $\frac{1}{p} + \frac{1}{q} = 1$. Thus

$$\begin{aligned}
 [\operatorname{tr}\{(A + B)^p\}]^{1/p} &= \max_{C \in \mathcal{S}} \operatorname{tr}\{(A + B)C\} \\
 &\leq \max_{C \in \mathcal{S}} \operatorname{tr}(AC) + \max_{C \in \mathcal{S}} \operatorname{tr}(BC) \\
 &= \{\operatorname{tr}(A^p)\}^{1/p} + \{\operatorname{tr}(B^p)\}^{1/p},
 \end{aligned}$$

with equality if and only if $\operatorname{tr}\{(A + B)C\}$, $\operatorname{tr}(AC)$, and $\operatorname{tr}(BC)$ are all maximized by the same C . According to Theorem 10.17, this requires that $(A + B)^p$, A^p , and B^p be proportional to one another or, equivalently, $B = cA$, for some scalar c . \square

10.6 THE ARITHMETIC-GEOMETRIC MEAN INEQUALITY

If $x_1 > 0$, $x_2 > 0$, and $0 \leq \alpha \leq 1$, then, since $g(x) = -\log(x)$ is a convex function on $x \in (0, \infty)$,

$$\begin{aligned}\log\{\alpha x_1 + (1 - \alpha)x_2\} &\geq \alpha \log(x_1) + (1 - \alpha) \log(x_2) \\ &= \log(x_1^\alpha x_2^{1-\alpha}),\end{aligned}$$

or, equivalently,

$$\alpha x_1 + (1 - \alpha)x_2 \geq x_1^\alpha x_2^{(1-\alpha)}.$$

Clearly this inequality also holds if either $x_1 = 0$ or $x_2 = 0$. More generally, the weighted arithmetic-geometric mean inequality is given by

$$\sum_{i=1}^m \alpha_i x_i \geq \prod_{i=1}^m x_i^{\alpha_i},$$

where $x_i \geq 0$, $\alpha_i \geq 0$, and $\sum_{i=1}^m \alpha_i = 1$. The special case in which $\alpha_i = \frac{1}{m}$ leads to

$$\frac{1}{m} \sum_{i=1}^m x_i \geq \left(\prod_{i=1}^m x_i \right)^{1/m},$$

which is known as the arithmetic-geometric mean inequality. We have equality in any of these inequalities if and only if $x_1 = \cdots = x_m$.

The first of our results in this section uses the arithmetic-geometric mean inequality to establish a relationship between the trace and determinant of a nonnegative definite matrix.

Theorem 10.21 Suppose A is an $m \times m$ nonnegative definite matrix. Then

$$\frac{1}{m} \operatorname{tr}(A) \geq |A|^{1/m},$$

with equality if and only if $A = cI_m$ for some $c > 0$.

Proof. Let $\lambda_1 \geq \cdots \geq \lambda_m$ be the eigenvalues of A . The inequality in Theorem 10.21 is a direct consequence of the arithmetic-geometric mean inequality since

$$\frac{1}{m} \operatorname{tr}(A) = \frac{1}{m} \sum_{i=1}^m \lambda_i \geq \left(\prod_{i=1}^m \lambda_i \right)^{1/m} = |A|^{1/m}.$$

We have equality if and only if $\lambda_1 = \cdots = \lambda_m$ or, equivalently, A is proportional to I_m . \square

The arithmetic-geometric mean inequality immediately extends to diagonal matrices. That is, if $\Lambda_1, \dots, \Lambda_n$ are $m \times m$ diagonal matrices with nonnegative diagonal elements, then

$$\frac{1}{n} \sum_{i=1}^n \Lambda_i - \left(\prod_{i=1}^n \Lambda_i \right)^{1/n}$$

is nonnegative definite and reduces to the null matrix if and only if $\Lambda_1 = \dots = \Lambda_n$. Our final result is a simple extension of this result to nonnegative definite matrices.

Theorem 10.22 Let A_1, \dots, A_n be $m \times m$ nonnegative definite matrices. If $A_i A_j = A_j A_i$ for all $i \neq j$, then the matrix

$$\frac{1}{n} \sum_{i=1}^n A_i - \left(\prod_{i=1}^n A_i \right)^{1/n}$$

is nonnegative definite and reduces to the null matrix if and only if $A_1 = \dots = A_n$.

Proof. By Theorem 4.19, there exists an orthogonal matrix P such that $A_i = P \Lambda_i P'$, where $\Lambda_i = \text{diag}(\lambda_{i1}, \dots, \lambda_{im})$ and $\lambda_{i1}, \dots, \lambda_{im}$ are the not necessarily ordered eigenvalues of A_i . Consequently,

$$\frac{1}{n} \sum_{i=1}^n A_i - \left(\prod_{i=1}^n A_i \right)^{1/n} = P \left\{ \frac{1}{n} \sum_{i=1}^n \Lambda_i - \left(\prod_{i=1}^n \Lambda_i \right)^{1/n} \right\} P',$$

and so the result follows directly from the result for diagonal matrices. \square

PROBLEMS

10.1 If $x \prec y$ and $y \prec x$, how are x and y related?

10.2 Suppose x , y , and z are all $m \times 1$ vectors, $x \prec y$, and the components of z are nonnegative.

(a) Show that

$$\sum_{i=1}^m y_{[i]} z_{[i]} \geq \sum_{i=1}^m x_{[i]} z_{[i]}.$$

(b) Show that

$$\sum_{i=1}^m y_{[i]} z_{[m+1-i]} \leq \sum_{i=1}^m x_{[i]} z_{[m+1-i]}.$$

10.3 Show that if $x \prec z$, $y \prec z$, and $0 \leq \lambda \leq 1$, then $\lambda x + (1 - \lambda)y \prec z$.

10.4 Suppose $x \prec y$. Show that there exists an orthogonal matrix Q such that $x = Py$, where $P = Q \odot Q$.

10.5 Show that $\mathbf{x} \prec \mathbf{y}$ if and only if

$$\sum_{i=1}^m |x_i - a| \leq \sum_{i=1}^m |y_i - a|$$

for all real scalars a .

10.6 Show that if for some $m \times m$ matrix P , $P\mathbf{x} \prec \mathbf{x}$ for every $m \times 1$ vector \mathbf{x} , then P is doubly stochastic.

10.7 If P is a nonsingular doubly stochastic matrix, show that $P^{-1}\mathbf{1}_m = \mathbf{1}_m$ and $\mathbf{1}'_m P^{-1} = \mathbf{1}'_m$. If, in addition, P^{-1} is doubly stochastic, that is, it has nonnegative elements, show that P is a permutation matrix.

10.8 Give an alternative proof of Theorem 10.3 by using the spectral decomposition of A .

10.9 Let A and B be $m \times m$ symmetric matrices having the partitioned forms

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A'_{12} & A_{22} \end{bmatrix}, \quad B = \begin{bmatrix} A_{11} & (0) \\ (0) & A_{22} \end{bmatrix},$$

where A_{11} is $m_1 \times m_1$, A_{22} is $m_2 \times m_2$ and $m_1 + m_2 = m$. If the components of \mathbf{a} are the eigenvalues of A and the components of \mathbf{b} are the eigenvalues of B , show that $\mathbf{b} \prec \mathbf{a}$.

10.10 Let \mathbf{x} be an $m \times 1$ vector with $x_i > 0$ for all i and $\sum_{i=1}^m x_i = 1$. Use Theorem 10.9 to show that

$$\sum_{i=1}^m \left(x_i + \frac{1}{x_i} \right)^a \geq \frac{(m^2 + 1)^a}{m^{a-1}}$$

for any $a > 1$.

10.11 Let \mathbf{x} be an $m \times 1$ vector with $x_i > 0$ for all i and $\bar{x} = m^{-1} \sum_{i=1}^m x_i$. Show that

$$\sum_{i=1}^m \frac{m\bar{x} - x_i}{x_i} \geq m(m-1).$$

10.12 Let \mathbf{x} and \mathbf{y} be $m \times 1$ vectors and A be an $m \times m$ nonnegative definite matrix. Show that

$$(\mathbf{x}'A\mathbf{y})^2 \leq (\mathbf{x}'A\mathbf{x})(\mathbf{y}'A\mathbf{y}),$$

with equality if and only if one of the vectors $A\mathbf{x}$ and $A\mathbf{y}$ is a scalar multiple of the other.

10.13 Let A be an $m \times m$ nonnegative definite matrix and \mathbf{x} and \mathbf{y} be $m \times 1$ vectors with \mathbf{y} being in the column space of A . Show that

$$(\mathbf{x}'\mathbf{y})^2 \leq (\mathbf{x}'A\mathbf{x})(\mathbf{y}'A^{-}\mathbf{y}),$$

with equality if and only if one of the vectors $A\mathbf{x}$ and \mathbf{y} is a scalar multiple of the other.

- 10.14** Let A and B be $m \times m$ positive definite matrices. Use Theorem 10.11 to show that

$$\mathbf{y}'(A+B)^{-1}\mathbf{y} \leq \frac{(\mathbf{y}'A^{-1}\mathbf{y})(\mathbf{y}'B^{-1}\mathbf{y})}{\mathbf{y}'(A^{-1}+B^{-1})\mathbf{y}}$$

for any $\mathbf{y} \neq \mathbf{0}$.

- 10.15** Let \mathbf{x} be an $m \times 1$ vector and A be an $m \times m$ positive definite matrix with eigenvalues $\lambda_1 \geq \cdots \geq \lambda_m$. Establish the Kantorovich inequality which is given by

$$\frac{(\mathbf{x}'A\mathbf{x})(\mathbf{x}'A^{-1}\mathbf{x})}{(\mathbf{x}'\mathbf{x})^2} \leq \frac{(\lambda_1 + \lambda_m)^2}{4\lambda_1\lambda_m}.$$

- 10.16** Let \mathbf{x} and \mathbf{y} be $m \times 1$ vectors. If $\mathbf{x}'\mathbf{1}_m = 1$ and $x_i \geq 0$ for $i = 1, \dots, m$, show that

$$(\mathbf{x}'\mathbf{y})^2 \leq \sum_{i=1}^m x_i y_i^2,$$

with equality if and only if $y_1 = \cdots = y_m$.

- 10.17** Let A be an $m \times m$ matrix with real eigenvalues. Show that

$$\{\operatorname{tr}(A)\}^2 \leq m \operatorname{tr}(A^2),$$

with equality if and only the eigenvalues are all equal.

- 10.18** If A and B are $m \times n$ matrices, show that

$$\operatorname{tr}\{(A'B)^2\} \leq \operatorname{tr}\{(A'A)(B'B)\},$$

with equality if and only if AB' is symmetric.

- 10.19** If A and B are $m \times n$ matrices, show that

$$\operatorname{tr}\{(A'B)^2\} \leq \operatorname{tr}\{(AA')(BB')\},$$

with equality if and only if $A'B$ is symmetric.

- 10.20** Let A be an $m \times m$ matrix. Show that

$$\operatorname{tr}(A^2) \leq \operatorname{tr}(A'A),$$

with equality if and only if A is symmetric.

- 10.21** Let A_1, \dots, A_n be $m \times m$ positive definite matrices, and let $\alpha_1, \dots, \alpha_n$ be nonnegative scalars satisfying $\sum_{i=1}^n \alpha_i = 1$. Show that

$$\prod_{i=1}^n |A_i|^{\alpha_i} \leq \left| \sum_{i=1}^n \alpha_i A_i \right|,$$

with equality if and only if $A_1 = \cdots = A_n$.

10.22 Let A and B be $m \times m$ nonnull nonnegative definite matrices and define

$$C = \alpha A + (1 - \alpha)B - A^\alpha B^{1-\alpha},$$

where $0 < \alpha < 1$. Use Theorem 10.18 to show that $\text{tr}(C) \geq 0$, with equality if and only if $A = B$, thereby establishing the inequality

$$\text{tr}(A^\alpha B^{1-\alpha}) \leq \text{tr}\{\alpha A + (1 - \alpha)B\}.$$

10.23 If A and B are $m \times m$ positive definite matrices and $0 < \alpha < 1$, show that

$$\text{tr}[\{\alpha A + (1 - \alpha)B\}^{-1}] \leq \alpha \text{tr}(A^{-1}) + (1 - \alpha)\text{tr}(B^{-1}).$$

10.24 Let A be an $m \times m$ nonnull nonnegative definite matrix. Show that

$$m^{-1} \text{tr}(AB) \geq |A|^{1/m}$$

for every positive definite $m \times m$ matrix B satisfying $|B| = 1$, with equality if and only if A is nonsingular and $B = |A|^{1/m} A^{-1}$.

10.25 Give an alternative proof of Theorem 10.19 by using the result from the previous problem.

10.26 Let A and B be $m \times m$ positive definite matrices and define

$$C = \begin{bmatrix} A & B \\ (0) & (0) \end{bmatrix}, \quad D = \begin{bmatrix} (0) & A'B \\ B'A & (0) \end{bmatrix}.$$

Compare the eigenvalues of $\frac{1}{2}C'C$ and D by using Theorem 3.28 to show that $\sigma_i(AB) \leq \sigma_i(A^2 + B^2)/2$ for $i = 1, \dots, m$, where $\sigma_1(A) \geq \dots \geq \sigma_m(A)$ denote the singular values of A , thereby extending the arithmetic-geometric mean inequality $ab \leq (a^2 + b^2)/2$.

11

SOME SPECIAL TOPICS RELATED TO QUADRATIC FORMS

11.1 INTRODUCTION

We have seen that if A is an $m \times m$ symmetric matrix and \mathbf{x} is an $m \times 1$ vector, then the function of \mathbf{x} , $\mathbf{x}'A\mathbf{x}$, is called a quadratic form in \mathbf{x} . In many statistical applications, \mathbf{x} is a random vector, whereas A is a matrix of constants. The most common situation is one in which \mathbf{x} has as its distribution, or as its asymptotic distribution, the multivariate normal distribution. In this chapter, we investigate some of the distributional properties of $\mathbf{x}'A\mathbf{x}$ in this setting. In particular, we are most interested in determining conditions under which $\mathbf{x}'A\mathbf{x}$ will have a chi-squared distribution.

11.2 SOME RESULTS ON IDEMPOTENT MATRICES

We have noted earlier that an $m \times m$ matrix A is said to be idempotent if $A^2 = A$. We will see in Section 11.3 that idempotent matrices play an essential role in the discussion of conditions under which a quadratic form in normal variates has a chi-squared distribution. Consequently, this section is devoted to establishing some of the basics results regarding idempotent matrices.

Theorem 11.1 Let A be an $m \times m$ idempotent matrix. Then

- (a) $I_m - A$ is also idempotent,

- (b) each eigenvalue of A is 0 or 1,
- (c) A is diagonalizable,
- (d) $\text{rank}(A) = \text{tr}(A)$.

Proof. Since $A^2 = A$, we have

$$(I_m - A)^2 = I_m - 2A + A^2 = I_m - A,$$

and so (a) holds. Let λ be an eigenvalue of A corresponding to the eigenvector \mathbf{x} , so that $A\mathbf{x} = \lambda\mathbf{x}$. Then because $A^2 = A$, we find that

$$\lambda\mathbf{x} = A\mathbf{x} = A^2\mathbf{x} = A(A\mathbf{x}) = A(\lambda\mathbf{x}) = \lambda A\mathbf{x} = \lambda^2\mathbf{x},$$

which implies that

$$\lambda(\lambda - 1)\mathbf{x} = \mathbf{0}.$$

Since eigenvectors are nonnull vectors, we must have $\lambda(\lambda - 1) = 0$, and so (b) follows. Let r be the number of eigenvalues of A equal to one, so that $m - r$ is the number of eigenvalues of A equal to zero. As a result, $A - I_m$ must have r eigenvalues equal to zero and $m - r$ eigenvalues equal to -1 . By Theorem 4.8, (c) will follow if we can show that

$$\text{rank}(A) = r, \quad \text{rank}(A - I_m) = m - r. \quad (11.1)$$

Now from Theorem 4.10, we know that the rank of any square matrix is at least as large as the number of its nonzero eigenvalues, so we must have

$$\text{rank}(A) \geq r, \quad \text{rank}(A - I_m) \geq m - r. \quad (11.2)$$

However, Corollary 2.10.1 gives

$$\begin{aligned} \text{rank}(A) + \text{rank}(I_m - A) &\leq \text{rank}\{A(I_m - A)\} + m \\ &= \text{rank}\{\mathbf{0}\} + m = m, \end{aligned}$$

which with (11.2) implies (11.1), so (c) is proven. Finally, (d) is an immediate consequence of (b) and (c). \square

Since any matrix with at least one 0 eigenvalue has to be a singular matrix, a nonsingular idempotent matrix has all of its eigenvalues equal to 1. However, the only diagonalizable matrix with all of its eigenvalues equal to 1 is the identity matrix; that is, the only nonsingular $m \times m$ idempotent matrix is I_m .

If A is a diagonal matrix, that is, A has the form $\text{diag}(a_1, \dots, a_m)$, then $A^2 = \text{diag}(a_1^2, \dots, a_m^2)$. Equating A and A^2 , we find that a diagonal matrix is idempotent

if and only if each diagonal element is 0 or 1, which is, of course, also an immediate consequence of Theorem 11.1(b).

Example 11.1 Although an idempotent matrix has each of its eigenvalues equal to 1 or 0, the converse is not true; that is, a matrix having only eigenvalues of 1 and 0 need not be an idempotent matrix. For instance, it is easily verified that the matrix

$$A = \begin{bmatrix} 0 & 1 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

has eigenvalues 0 and 1 with multiplicities 2 and 1, respectively. However,

$$A^2 = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

so that A is not idempotent.

The matrix A in Example 11.1 is not idempotent because it is not diagonalizable. In other words, an $m \times m$ matrix A is idempotent if and only if each of its eigenvalues is 0 or 1 and it is diagonalizable. In fact, we have a special case in Theorem 11.2.

Theorem 11.2 Let A be an $m \times m$ symmetric matrix. Then A is idempotent if and only if each eigenvalue of A is 0 or 1.

Proof. Let $A = X\Lambda X'$ be the spectral decomposition of A , so that X is an orthogonal matrix and Λ is diagonal. Then

$$A^2 = (X\Lambda X')^2 = X\Lambda X'X\Lambda X' = X\Lambda^2 X'.$$

Clearly, this equals A if and only if each diagonal element of Λ , that is, each eigenvalue of A , is 0 or 1. \square

Theorem 11.2 is generalized in Corollary 11.2.1.

Corollary 11.2.1 Let A be an $m \times m$ symmetric matrix and B be an $m \times m$ positive definite matrix. Then AB is idempotent if and only if each eigenvalue of AB is 0 or 1.

Proof. Since B is positive definite, it can be expressed as $B = TT'$, where T is an $m \times m$ nonsingular matrix. Note that if the equation

$$ABAB = AB \tag{11.3}$$

is premultiplied by T' and postmultiplied by $T^{-1'}$, it yields

$$T'ATT'AT = T'AT. \quad (11.4)$$

Conversely, premultiplying (11.4) by $T^{-1'}$ and postmultiplying by T' , we get (11.3). That is, AB is idempotent if and only if $T'AT$ is idempotent. Since AB and $T'AT$ have the same eigenvalues by Theorem 3.2(d), the result follows immediately from Theorem 11.2. \square

Theorem 11.3 gives some conditions for the sum of two idempotent matrices and the product of two idempotent matrices to be idempotent.

Theorem 11.3 Let A and B be $m \times m$ idempotent matrices. Then

- (a) $A + B$ is idempotent if and only if $AB = BA = (0)$,
- (b) AB is idempotent if $AB = BA$.

Proof. Since A and B are idempotent, we have

$$(A + B)^2 = A^2 + B^2 + AB + BA = A + B + AB + BA,$$

so that $A + B$ will be idempotent if and only if

$$AB = -BA. \quad (11.5)$$

Premultiplication of (11.5) by B and postmultiplication by A yields the identity

$$(BA)^2 = -BA, \quad (11.6)$$

since A and B are idempotent. Similarly, premultiplying (11.5) by A and postmultiplying by B , we also find that

$$(AB)^2 = -AB. \quad (11.7)$$

Thus, it follows from (11.6) and (11.7) that both $-BA$ and $-AB$ are idempotent matrices, and because of (11.5), so then are AB and BA . Part (a) now follows because the null matrix is the only idempotent matrix whose negative is also idempotent. To prove (b), note that if A and B commute under multiplication, then

$$(AB)^2 = ABAB = A(BA)B = A(AB)B = A^2B^2 = AB,$$

and so the result follows. \square

Example 11.2 The conditions given in Theorem 11.3 for $(A + B)$ to be idempotent are necessary and sufficient, whereas the condition given for AB to be idempotent is

only sufficient. We can illustrate that this second condition is not necessary through a simple example. Let A and B be defined as

$$A = \begin{bmatrix} 1 & 1 \\ 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 \\ 1 & 1 \end{bmatrix},$$

and observe that $A^2 = A$ and $B^2 = B$, so that both A and B are idempotent. In addition, $AB = A$, so that AB is also idempotent. However, $AB \neq BA$ because $BA = B$.

It is easily verified that a matrix having the form $A = B(C'B)^{-1}C'$ is idempotent. Our next result shows that every idempotent matrix can be expressed in this way.

Theorem 11.4 Suppose A is an $m \times m$ idempotent matrix of rank r . Then there exist $m \times r$ matrices B and C such that

$$A = B(C'B)^{-1}C'.$$

Proof. It follows from the singular value decomposition that A can be expressed as $A = BDC'$, where B and C are $m \times r$ full-rank matrices and D is a nonsingular $r \times r$ matrix. Since A is idempotent, we have

$$BDC' = BDC'BDC'.$$

Premultiplying this by $(B'B)^{-1}B'$ and postmultiplying by $C(C'C)^{-1}$ yields

$$D = DC'BD,$$

from which we get $D = (C'B)^{-1}$ as required. \square

Most of the statistical applications involving idempotent matrices deal with symmetric idempotent matrices. For this reason, we end this section with some results for this special class of matrices. Theorem 11.5 gives some restrictions on the elements of a symmetric idempotent matrix.

Theorem 11.5 Suppose A is an $m \times m$ symmetric idempotent matrix. Then

- (a) $a_{ii} \geq 0$ for $i = 1, \dots, m$,
- (b) $a_{ii} \leq 1$ for $i = 1, \dots, m$,
- (c) $a_{ij} = a_{ji} = 0$, for all $j \neq i$, if $a_{ii} = 0$ or $a_{ii} = 1$.

Proof. Since A is idempotent and symmetric, it follows that

$$\begin{aligned} a_{ii} &= (A)_{ii} = (A^2)_{ii} = (A'A)_{ii} \\ &= (A')_{i \cdot} (A)_{\cdot i} = \sum_{j=1}^m a_{ji}^2, \end{aligned} \tag{11.8}$$

which clearly must be nonnegative. In addition, from (11.8), we have

$$a_{ii} = a_{ii}^2 + \sum_{j \neq i} a_{ji}^2,$$

so that $a_{ii} \geq a_{ii}^2$ or $a_{ii}(1 - a_{ii}) \geq 0$. However, because a_{ii} is nonnegative, this leads to $(1 - a_{ii}) \geq 0$, and thus (b) must hold. If $a_{ii} = 0$ or $a_{ii} = 1$, then $a_{ii} = a_{ii}^2$, and so again using (11.8), we must have

$$\sum_{j \neq i} a_{ji}^2 = 0,$$

which, along with the symmetry of A , establishes (c). \square

Theorem 11.6 is useful in those situations in which it is easier to verify an identity such as $A^3 = A^2$ than the identity $A^2 = A$.

Theorem 11.6 Suppose that for some positive integer i , the $m \times m$ symmetric matrix A satisfies $A^{i+1} = A^i$. Then A is an idempotent matrix.

Proof. If $\lambda_1, \dots, \lambda_m$ are the eigenvalues of A , then $\lambda_1^{i+1}, \dots, \lambda_m^{i+1}$ and $\lambda_1^i, \dots, \lambda_m^i$ are the eigenvalues of A^{i+1} and A^i , respectively. However, the identity $A^{i+1} = A^i$ implies that $\lambda_j^{i+1} = \lambda_j^i$, for $j = 1, \dots, m$, so each λ_j must be either 0 or 1. The result now follows from Theorem 11.2. \square

11.3 COCHRAN'S THEOREM

Theorem 11.7, sometimes referred to as Cochran's Theorem (Cochran, 1934), will be useful in establishing the independence of several different quadratic forms in the same normal variables.

Theorem 11.7 Let each of the $m \times m$ matrices A_1, \dots, A_k be symmetric and idempotent, and suppose that $A_1 + \dots + A_k = I_m$. Then $A_i A_j = (0)$ whenever $i \neq j$.

Proof. Select any one of the matrices, say A_h , and denote its rank by r . Since A_h is symmetric and idempotent, an orthogonal matrix P exists, such that

$$P' A_h P = \text{diag}(I_r, (0)).$$

For $j \neq h$, define $B_j = P' A_j P$, and note that

$$\begin{aligned} I_m &= P' I_m P = P' \left(\sum_{j=1}^k A_j \right) P = \left(\sum_{j=1}^k P' A_j P \right) \\ &= \text{diag}(I_r, (0)) + \sum_{j \neq h} B_j, \end{aligned}$$

or equivalently,

$$\sum_{j \neq h} B_j = \text{diag}((0), I_{m-r}).$$

In particular, for $l = 1, \dots, r$,

$$\sum_{j \neq h} (B_j)_{ll} = 0.$$

However, clearly B_j is symmetric and idempotent because A_j is, and so, from Theorem 11.5(a), its diagonal elements are nonnegative. Thus, we must have $(B_j)_{ll} = 0$ for each $l = 1, \dots, r$, which, along with Theorem 11.5(c), implies that B_j must be of the form

$$B_j = \text{diag}((0), C_j),$$

where C_j is an $(m-r) \times (m-r)$ symmetric idempotent matrix. Now, for any $j \neq h$,

$$\begin{aligned} P' A_h A_j P &= (P' A_h P)(P' A_j P) \\ &= \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} (0) & (0) \\ (0) & C_j \end{bmatrix} = (0), \end{aligned}$$

which can be true only if $A_h A_j = (0)$, because P is nonsingular. Our proof is now complete, because h was arbitrary. \square

Our next result is an extension of Cochran's Theorem.

Theorem 11.8 Let A_1, \dots, A_k be $m \times m$ symmetric matrices, and define $A = A_1 + \dots + A_k$. Consider the following statements:

- (a) A_i is idempotent for $i = 1, \dots, k$.
- (b) A is idempotent.
- (c) $A_i A_j = (0)$, for all $i \neq j$.

Then if any two of these conditions hold, the third condition must also hold.

Proof. First we show that (a) and (b) imply (c). Since A is symmetric and idempotent, an orthogonal matrix P exists, such that

$$P' A P = P' (A_1 + \dots + A_k) P = \text{diag}(I_r, (0)), \quad (11.9)$$

where $r = \text{rank}(A)$. Let $B_i = P' A_i P$ for $i = 1, \dots, k$, and note that B_i is symmetric and idempotent. Thus, it follows from (11.9) and Theorem 11.5 that B_i must be

of the form $\text{diag}(C_i, (0))$, where the $r \times r$ matrix C_i also must be symmetric and idempotent. However, (11.9) also implies that

$$C_1 + \cdots + C_k = I_r.$$

Consequently, C_1, \dots, C_k satisfy the conditions of Theorem 11.7, and so $C_i C_j = (0)$ for every $i \neq j$. From this result, we get $B_i B_j = (0)$ and, hence, $A_i A_j = (0)$ for every $i \neq j$, as is required. That (a) and (c) imply (b) follows immediately, because

$$\begin{aligned} A^2 &= \left(\sum_{i=1}^k A_i \right)^2 = \sum_{i=1}^k \sum_{j=1}^k A_i A_j = \sum_{i=1}^k A_i^2 + \sum_{i \neq j} A_i A_j \\ &= \sum_{i=1}^k A_i = A. \end{aligned}$$

Finally, we must prove that (b) and (c) imply (a). If (c) holds, then $A_i A_j = A_j A_i$ for all $i \neq j$, and so by Theorem 4.19, the matrices A_1, \dots, A_k can be simultaneously diagonalized; that is, an orthogonal matrix Q exists, such that

$$Q' A_i Q = D_i,$$

where each of the matrices D_1, \dots, D_k is diagonal. Furthermore,

$$D_i D_j = Q' A_i Q Q' A_j Q = Q' A_i A_j Q = Q' (0) Q = (0), \quad (11.10)$$

for every $i \neq j$. Now because A is symmetric and idempotent, so also is the diagonal matrix

$$Q' A Q = D_1 + \cdots + D_k.$$

As a result, each diagonal element of $Q' A Q$ must be either 0 or 1, and because of (11.10), the same can be said of the diagonal elements of D_1, \dots, D_k . Thus, for each i , D_i is symmetric and idempotent and, hence, so is $A_i = Q D_i Q'$. This completes the proof. \square

Suppose that the three conditions given in Theorem 11.8 hold. Then (a) implies that $\text{tr}(A_i) = \text{rank}(A_i)$ and (b) implies that

$$\text{rank}(A) = \text{tr}(A) = \text{tr} \left(\sum_{i=1}^k A_i \right) = \sum_{i=1}^k \text{tr}(A_i) = \sum_{i=1}^k \text{rank}(A_i).$$

Thus, we have shown that the conditions in Theorem 11.8 imply the fourth condition

$$(d) \text{ rank}(A) = \sum_{i=1}^k \text{rank}(A_i).$$

Conversely, suppose that conditions (b) and (d) hold. We will show that these imply (a) and (c). Let $H = \text{diag}(A_1, \dots, A_k)$ and $F = \mathbf{1}_m \otimes I_m$, so that $A = F'HF$. Then (d) can be written $\text{rank}(F'HF) = \text{rank}(H)$, and so it follows from Theorem 5.26 that $F(F'HF)^-F'$ is a generalized inverse of H for any generalized inverse $(F'HF)^-$ of $F'HF$. However, because A is idempotent, $AI_mA = A$, and hence, I_m is a generalized inverse of $A = F'HF$. Thus, FF' is a generalized inverse of H , which yields the equation

$$HFF'H = H,$$

which in partitioned form is

$$\begin{bmatrix} A_1^2 & A_1A_2 & \cdots & A_1A_k \\ A_2A_1 & A_2^2 & \cdots & A_2A_k \\ \vdots & \vdots & & \vdots \\ A_kA_1 & A_kA_2 & \cdots & A_k^2 \end{bmatrix} = \begin{bmatrix} A_1 & (0) & \cdots & (0) \\ (0) & A_2 & \cdots & (0) \\ \vdots & \vdots & & \vdots \\ (0) & (0) & \cdots & A_k \end{bmatrix}.$$

This equation immediately gives conditions (a) and (c). Corollary 11.8.1 summarizes the relationship among these four conditions.

Corollary 11.8.1 Let A_1, \dots, A_k be $m \times m$ symmetric matrices, and define $A = A_1 + \dots + A_k$. Consider the following statements.

- (a) A_i is idempotent for $i = 1, \dots, k$.
- (b) A is idempotent.
- (c) $A_iA_j = (0)$, for all $i \neq j$.
- (d) $\text{rank}(A) = \sum_{i=1}^k \text{rank}(A_i)$.

All four of the conditions hold if any two of (a), (b), and (c) hold, or if (b) and (d) hold.

11.4 DISTRIBUTION OF QUADRATIC FORMS IN NORMAL VARIATES

The relationship between the normal and chi-squared distributions is fundamental in obtaining the distribution of a quadratic form in normal random variables. Recall that if z_1, \dots, z_r are independent random variables with $z_i \sim N(0, 1)$ for each i , then

$$\sum_{i=1}^r z_i^2 \sim \chi_r^2.$$

This is used in Theorem 11.9 to determine when the quadratic form $\mathbf{x}'A\mathbf{x}$ has a chi-squared distribution if the components of \mathbf{x} are independently distributed, each having the $N(0, 1)$ distribution.

Theorem 11.9 Let $\mathbf{x} \sim N_m(\mathbf{0}, I_m)$, and suppose that the $m \times m$ matrix A is symmetric, is idempotent, and has rank r . Then $\mathbf{x}'A\mathbf{x} \sim \chi_r^2$.

Proof. Since A is symmetric and idempotent, an orthogonal matrix P exists, such that

$$A = PDP',$$

where $D = \text{diag}(I_r, (0))$. Let $\mathbf{z} = P'\mathbf{x}$, and note that because $\mathbf{x} \sim N_m(\mathbf{0}, I_m)$,

$$E(\mathbf{z}) = E(P'\mathbf{x}) = P'E(\mathbf{x}) = P'\mathbf{0} = \mathbf{0},$$

$$\text{var}(\mathbf{z}) = \text{var}(P'\mathbf{x}) = P'\{\text{var}(\mathbf{x})\}P = P'I_mP = P'P = I_m,$$

and so $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$; that is, the components of \mathbf{z} are, like the components of \mathbf{x} , independent standard normal random variables. Now because of the form of D , we find that

$$\mathbf{x}'A\mathbf{x} = \mathbf{x}'PDP'\mathbf{x} = \mathbf{z}'D\mathbf{z} = \sum_{i=1}^r z_i^2,$$

and so the result then follows. \square

Theorem 11.9 is a special case of Theorem 11.10 in which the multivariate normal distribution has a general nonsingular covariance matrix.

Theorem 11.10 Let $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is a positive definite matrix, and let A be an $m \times m$ symmetric matrix. If $A\Omega$ is idempotent and $\text{rank}(A\Omega) = r$, then $\mathbf{x}'A\mathbf{x} \sim \chi_r^2$.

Proof. Since Ω is positive definite, a nonsingular matrix T exists, which satisfies $\Omega = TT'$. If we define $\mathbf{z} = T^{-1}\mathbf{x}$, then $E(\mathbf{z}) = T^{-1}E(\mathbf{x}) = \mathbf{0}$ and

$$\begin{aligned} \text{var}(\mathbf{z}) &= \text{var}(T^{-1}\mathbf{x}) = T^{-1}\{\text{var}(\mathbf{x})\}T^{-1'} \\ &= T^{-1}(TT')T^{-1'} = I_m, \end{aligned}$$

so that $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$. The quadratic form $\mathbf{x}'A\mathbf{x}$ can be written in terms of \mathbf{z} because

$$\mathbf{x}'A\mathbf{x} = \mathbf{x}'T^{-1'}T'ATT^{-1}\mathbf{x} = \mathbf{z}'T'AT\mathbf{z}.$$

All that remains is to show that $T'AT$ satisfies the conditions of Theorem 11.9. Clearly, $T'AT$ is symmetric, because A is, and idempotent because

$$(T'AT)^2 = T'ATT'AT = T'A\Omega AT = T'AT,$$

where the last equality follows from the identity $A\Omega A = A$, which is a consequence of the fact that $A\Omega$ is idempotent and Ω is nonsingular. Finally, because $T'AT$ and

$A\Omega$ are idempotent, we have

$$\begin{aligned}\text{rank}(T'AT) &= \text{tr}(T'AT) = \text{tr}(ATT') \\ &= \text{tr}(A\Omega) = \text{rank}(A\Omega) = r,\end{aligned}$$

and so the proof is complete. \square

It is not uncommon to have a quadratic form in a vector that has a singular multivariate normal distribution. Our next result generalizes Theorem 11.10 to this situation.

Theorem 11.11 Let $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is positive semidefinite, and suppose that A is an $m \times m$ symmetric matrix. If $\Omega A \Omega A \Omega = \Omega A \Omega$ and $\text{tr}(A\Omega) = r$, then $\mathbf{x}'A\mathbf{x} \sim \chi_r^2$.

Proof. Let $n = \text{rank}(\Omega)$, where $n < m$. Then an $m \times m$ orthogonal matrix $P = [P_1 \ P_2]$ exists, such that

$$\Omega = [P_1 \ P_2] \begin{bmatrix} \Lambda & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} P_1' \\ P_2' \end{bmatrix} = P_1 \Lambda P_1',$$

where P_1 is $m \times n$ and Λ is an $n \times n$ nonsingular diagonal matrix. Define

$$\mathbf{z} = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} P_1' \mathbf{x} \\ P_2' \mathbf{x} \end{bmatrix} = P' \mathbf{x},$$

and note that because $P'\mathbf{0} = \mathbf{0}$ and $P'\Omega P = \text{diag}(\Lambda, (0))$, $\mathbf{z} \sim N_m(\mathbf{0}, \text{diag}(\Lambda, (0)))$, which means that $\mathbf{z} = (z_1', \mathbf{0}')'$, where z_1 has the nonsingular distribution $N_n(\mathbf{0}, \Lambda)$. Now

$$\mathbf{x}'A\mathbf{x} = \mathbf{x}'PP'APP'\mathbf{x} = \mathbf{z}'P'AP\mathbf{z} = \mathbf{z}_1'P_1'AP_1\mathbf{z}_1,$$

and so the proof will be complete if we can show that the symmetric matrix $P_1'AP_1$ satisfies the conditions of the previous theorem, namely, that $P_1'AP_1\Lambda$ is idempotent and $\text{rank}(P_1'AP_1\Lambda) = r$. Since $\Omega A \Omega A \Omega = \Omega A \Omega$, we have

$$\begin{aligned}(\Lambda^{1/2}P_1'AP_1\Lambda^{1/2})^3 &= \Lambda^{1/2}P_1'A(P_1\Lambda P_1')A(P_1\Lambda P_1')AP_1\Lambda^{1/2} \\ &= \Lambda^{1/2}P_1'A\Omega A\Omega AP_1\Lambda^{1/2} = \Lambda^{1/2}P_1'A\Omega AP_1\Lambda^{1/2} \\ &= \Lambda^{1/2}P_1'A(P_1\Lambda P_1')AP_1\Lambda^{1/2} = (\Lambda^{1/2}P_1'AP_1\Lambda^{1/2})^2,\end{aligned}$$

and so the idempotency of $\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}$ follows from Theorem 11.6. However, this also establishes the idempotency of $P_1'AP_1\Lambda$ because Λ is nonsingular. Its rank is r because

$$\text{rank}(P_1'AP_1\Lambda) = \text{tr}(P_1'AP_1\Lambda) = \text{tr}(AP_1\Lambda P_1') = \text{tr}(A\Omega) = r,$$

and so the result follows. \square

Until now, our results have dealt with normal distributions having the zero mean vector. In some applications, such as the determination of nonnull distributions in hypothesis testing situations, we encounter quadratic forms in normal random vectors having nonzero means. The next two theorems are helpful in determining whether such a quadratic form has a chi-squared distribution. The proof of the first of these two theorems, which is very similar to that of Theorem 11.10, is left to the reader. It applies the relationship between the normal distribution and the noncentral chi-squared distribution; that is, if y_1, \dots, y_r are independently distributed with $y_i \sim N(\mu_i, 1)$, then

$$\sum_{i=1}^r y_i^2 \sim \chi_r^2(\lambda),$$

where the noncentrality parameter of this noncentral chi-squared distribution is given by

$$\lambda = \frac{1}{2} \sum_{i=1}^r \mu_i^2.$$

Theorem 11.12 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is a positive definite matrix, and let A be an $m \times m$ symmetric matrix. If $A\Omega$ is idempotent and $\text{rank}(A\Omega) = r$, then $\mathbf{x}'A\mathbf{x} \sim \chi_r^2(\lambda)$, where $\lambda = \frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu}$.

Theorem 11.13 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive semidefinite of rank n , and suppose that A is an $m \times m$ symmetric matrix. Then $\mathbf{x}'A\mathbf{x} \sim \chi_r^2(\lambda)$, where $\lambda = \frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu}$ if

- (a) $\Omega A \Omega A \Omega = \Omega A \Omega$,
- (b) $\boldsymbol{\mu}' A \Omega A \Omega = \boldsymbol{\mu}' A \Omega$,
- (c) $\boldsymbol{\mu}' A \Omega A \boldsymbol{\mu} = \boldsymbol{\mu}' A \boldsymbol{\mu}$,
- (d) $\text{tr}(A\Omega) = r$.

Proof. Let P_1, P_2 , and Λ be defined as in the proof of Theorem 11.11, so that $\Omega = P_1 \Lambda P_1'$. Put $C = [P_1 \Lambda^{-1/2} \quad P_2]$, and note that

$$\begin{aligned} \mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} &= \begin{bmatrix} \Lambda^{-1/2} P_1' \mathbf{x} \\ P_2' \mathbf{x} \end{bmatrix} \\ &= C' \mathbf{x} \sim N_m \left(\begin{bmatrix} \Lambda^{-1/2} P_1' \boldsymbol{\mu} \\ P_2' \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} I_n & (0) \\ (0) & (0) \end{bmatrix} \right). \end{aligned}$$

In other words,

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ P_2' \boldsymbol{\mu} \end{bmatrix},$$

where $\mathbf{z}_1 \sim N_n(\Lambda^{-1/2}P_1'\boldsymbol{\mu}, I_n)$. Now because $C^{-1'} = [P_1\Lambda^{1/2} \quad P_2]$, we find that

$$\begin{aligned} \mathbf{x}'\mathbf{A}\mathbf{x} &= \mathbf{x}'CC^{-1}AC^{-1'}C'\mathbf{x} = \mathbf{z}'C^{-1}AC^{-1'}\mathbf{z} \\ &= [\mathbf{z}_1' \quad \boldsymbol{\mu}'P_2] \begin{bmatrix} \Lambda^{1/2}P_1'AP_1\Lambda^{1/2} & \Lambda^{1/2}P_1'AP_2 \\ P_2'AP_1\Lambda^{1/2} & P_2'AP_2 \end{bmatrix} \begin{bmatrix} \mathbf{z}_1 \\ P_2'\boldsymbol{\mu} \end{bmatrix} \\ &= \mathbf{z}_1'\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}\mathbf{z}_1 + \boldsymbol{\mu}'P_2P_2'AP_2P_2'\boldsymbol{\mu} \\ &\quad + 2\boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}\mathbf{z}_1. \end{aligned} \quad (11.11)$$

However, conditions (a)–(c) imply the identities

- (i) $P_1'AP_1 = P_1'AP_1$,
- (ii) $\boldsymbol{\mu}'P_2P_2'AP_1 = \boldsymbol{\mu}'P_2P_2'AP_1$,
- (iii) $\boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}\mathbf{z}_1 = \boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}\mathbf{z}_1$,

in particular, (a) implies (i), (b) and (i) imply (ii), whereas (iii) follows from (b), (c), (i), and (ii). Using these identities in (11.11), we obtain

$$\begin{aligned} \mathbf{x}'\mathbf{A}\mathbf{x} &= \mathbf{z}_1'\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}\mathbf{z}_1 + \boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}\mathbf{z}_1 \\ &\quad + 2\boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}\mathbf{z}_1 \\ &= (\mathbf{z}_1 + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu})'\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}(\mathbf{z}_1 + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu}) \\ &= \mathbf{w}'A_*\mathbf{w}. \end{aligned}$$

Now, $\mathbf{w} = (\mathbf{z}_1 + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu}) \sim N_n(\boldsymbol{\theta}, I_n)$, where

$$\boldsymbol{\theta} = \Lambda^{-1/2}P_1'\boldsymbol{\mu} + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu},$$

and, because $A_* = \Lambda^{1/2}P_1'AP_1\Lambda^{1/2}$ is idempotent, a consequence of (i), we may apply Theorem 11.12; that is, $\mathbf{w}'A_*\mathbf{w} \sim \chi_r^2(\lambda)$, where

$$r = \text{tr}(A_*I_n) = \text{tr}(\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}) = \text{tr}(AP_1\Lambda P_1') = \text{tr}(A\Omega)$$

and

$$\begin{aligned} \lambda &= \frac{1}{2}\boldsymbol{\theta}'A_*\boldsymbol{\theta} = \frac{1}{2}(\Lambda^{-1/2}P_1'\boldsymbol{\mu} + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu})' \\ &\quad \times \Lambda^{1/2}P_1'AP_1\Lambda^{1/2}(\Lambda^{-1/2}P_1'\boldsymbol{\mu} + \Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu}) \\ &= \frac{1}{2}(\boldsymbol{\mu}'P_1P_1'AP_1P_1'\boldsymbol{\mu} + \boldsymbol{\mu}'P_2P_2'AP_1\Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu} + 2\boldsymbol{\mu}'P_1P_1'AP_1\Lambda^{1/2}P_1'AP_2P_2'\boldsymbol{\mu}) \\ &= \frac{1}{2}(\boldsymbol{\mu}'P_1P_1'AP_1P_1'\boldsymbol{\mu} + \boldsymbol{\mu}'P_2P_2'AP_2P_2'\boldsymbol{\mu} + 2\boldsymbol{\mu}'P_1P_1'AP_2P_2'\boldsymbol{\mu}) \\ &= \frac{1}{2}\boldsymbol{\mu}'(P_1P_1' + P_2P_2')A(P_1P_1' + P_2P_2')\boldsymbol{\mu} = \frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu}. \end{aligned}$$

This completes the proof. \square

A matrix A satisfying conditions (a), (b), and (c) of Theorem 11.13 is Ω^+ , the Moore–Penrose inverse of Ω . That is, if $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, then $\mathbf{x}'\Omega^+\mathbf{x}$ will have a chi-squared distribution because the identity $\Omega^+\Omega\Omega^+ = \Omega^+$ ensures that conditions (a), (b), and (c) hold. The degrees of freedom $r = \text{rank}(\Omega)$ because $\text{rank}(\Omega^+\Omega) = \text{rank}(\Omega)$.

All of the theorems presented in this section give sufficient conditions for a quadratic form to have a chi-squared distribution. Actually, in each case, the stated conditions are necessary conditions as well, which is most easily proven using moment generating functions. For more details, see Mathai and Provost (1992) or Searle (1971).

Example 11.3 Let x_1, \dots, x_n be a random sample from a normal distribution with mean μ and variance σ^2 ; that is, the x_i 's are independent random variables, each having the distribution $N(\mu, \sigma^2)$. The sample variance s^2 is given by

$$s^2 = \frac{1}{(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2.$$

We will use the results of this section to show that

$$t = \frac{(n-1)s^2}{\sigma^2} = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} \sim \chi_{n-1}^2.$$

Define the $n \times 1$ vector $\mathbf{x} = (x_1, \dots, x_n)'$, so that $\mathbf{x} \sim N_n(\mu \mathbf{1}_n, \sigma^2 I_n)$. Note that if the $n \times n$ matrix $A = (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') / \sigma^2$, then

$$\begin{aligned} \mathbf{x}' A \mathbf{x} &= \frac{\{\mathbf{x}' \mathbf{x} - n^{-1} (\mathbf{1}_n' \mathbf{x})^2\}}{\sigma^2} = \sigma^{-2} \left\{ \sum_{i=1}^n x_i^2 - n^{-1} \left(\sum_{i=1}^n x_i \right)^2 \right\} \\ &= \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1)s^2}{\sigma^2} = t, \end{aligned}$$

and so t is a quadratic form in the random vector \mathbf{x} . The matrix $A(\sigma^2 I_n) = \sigma^2 A$ is idempotent because

$$\begin{aligned} (\sigma^2 A)^2 &= (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')^2 = I_n - 2n^{-1} \mathbf{1}_n \mathbf{1}_n' + n^{-2} \mathbf{1}_n \mathbf{1}_n' \mathbf{1}_n \mathbf{1}_n' \\ &= I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n' = \sigma^2 A, \end{aligned}$$

and so, by Theorem 11.12, t has a chi-squared distribution. This chi-squared distribution has $n-1$ degrees of freedom because

$$\begin{aligned} \text{tr}(\sigma^2 A) &= \text{tr}(I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') = \text{tr}(I_n) - n^{-1} \text{tr}(\mathbf{1}_n \mathbf{1}_n') \\ &= n - n^{-1} \mathbf{1}_n' \mathbf{1}_n = n - 1, \end{aligned}$$

and the noncentrality parameter is given by

$$\begin{aligned}\lambda &= \frac{1}{2} \boldsymbol{\mu}' A \boldsymbol{\mu} = \frac{1}{2} \frac{\mu^2}{\sigma^2} \mathbf{1}'_n (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) \mathbf{1}_n \\ &= \frac{1}{2} \frac{\mu^2}{\sigma^2} (\mathbf{1}'_n \mathbf{1}_n - n^{-1} \mathbf{1}'_n \mathbf{1}_n \mathbf{1}'_n \mathbf{1}_n) \\ &= \frac{1}{2} \frac{\mu^2}{\sigma^2} (n - n) = 0.\end{aligned}$$

Thus, we have shown that $t \sim \chi^2_{n-1}$.

11.5 INDEPENDENCE OF QUADRATIC FORMS

We now consider the situation in which we have several different quadratic forms, each a function of the same multivariate normal vector. In some settings, it is important to be able to determine whether these quadratic forms are distributed independently of one another. For instance, this is useful in the partitioning of chi-squared random variables as well as in the formation of ratios having an F distribution.

We begin with the following basic result regarding the statistical independence of two quadratic forms in the same normal vector.

Theorem 11.14 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite, and suppose that A and B are $m \times m$ symmetric matrices. If $A\Omega B = (0)$, then $\mathbf{x}'A\mathbf{x}$ and $\mathbf{x}'B\mathbf{x}$ are independently distributed.

Proof. Since Ω is positive definite, a nonsingular matrix T exists, such that $\Omega = TT'$. Define $G = T'AT$ and $H = T'BT$, and note that if $A\Omega B = (0)$, then

$$GH = (T'AT)(T'BT) = T'A\Omega BT = T'(0)T = (0). \quad (11.12)$$

Consequently, because of the symmetry of G and H , we also have

$$(0) = (0)' = (GH)' = H'G' = HG,$$

and so we have established that $GH = HG$. From Theorem 4.18, we know that an orthogonal matrix P exists that simultaneously diagonalizes G and H ; that is, for some diagonal matrices C and D ,

$$P'GP = P'T'ATP = C, \quad P'HP = P'T'BT = D. \quad (11.13)$$

However, using (11.12) and (11.13), we find that

$$(0) = GH = PCP'PDP' = PCDP',$$

which can be true only if $CD = (0)$. Since C and D are diagonal matrices, this means that if the i th diagonal element of one of these matrices is nonzero, the i th diagonal element of the other must be zero. As a result, by choosing P appropriately, we may obtain C and D in the form $C = \text{diag}(c_1, \dots, c_{m_1}, 0, \dots, 0)$ and $D = \text{diag}(0, \dots, 0, d_{m_1+1}, \dots, d_m)$ for some integer m_1 . If we let $\mathbf{y} = P'T^{-1}\mathbf{x}$, then our two quadratic forms simplify as

$$\mathbf{x}'A\mathbf{x} = \mathbf{x}'T^{-1'}PP'T'ATPP'T^{-1}\mathbf{x} = \mathbf{y}'C\mathbf{y} = \sum_{i=1}^{m_1} c_i y_i^2$$

and

$$\mathbf{x}'B\mathbf{x} = \mathbf{x}'T^{-1'}PP'T'BTTPP'T^{-1}\mathbf{x} = \mathbf{y}'D\mathbf{y} = \sum_{i=m_1+1}^m d_i y_i^2;$$

that is, the first quadratic form is a function only of y_1, \dots, y_{m_1} , whereas the second quadratic form is a function of y_{m_1+1}, \dots, y_m . Since

$$\text{var}(\mathbf{y}) = \text{var}(P'T^{-1}\mathbf{x}) = P'T^{-1}\Omega T^{-1'}P = I_m,$$

the result then follows from the independence of y_1, \dots, y_m , which is a consequence of the fact that \mathbf{y} is normal. \square

Example 11.4 Suppose that $\mathbf{x}_1, \dots, \mathbf{x}_k$ are independently distributed with $\mathbf{x}_i = (x_{i1}, \dots, x_{in})' \sim N_n(\mu \mathbf{1}_n, \sigma^2 I_n)$ for each i . Let t_1 and t_2 be the random quantities defined by

$$t_1 = n \sum_{i=1}^k (\bar{x}_i - \bar{x})^2, \quad t_2 = \sum_{i=1}^k \sum_{j=1}^n (x_{ij} - \bar{x}_i)^2,$$

where

$$\bar{x}_i = \sum_{j=1}^n \frac{x_{ij}}{n}, \quad \bar{x} = \sum_{i=1}^k \frac{\bar{x}_i}{k}.$$

Note that t_1 and t_2 are the formulas for the sum of squares for treatments and the sum of squares for error in a balanced one-way classification model (Example 8.2). Now t_1 can be expressed as

$$t_1 = n \left\{ \sum_{i=1}^k \bar{x}_i^2 - k^{-1} \left(\sum_{i=1}^k \bar{x}_i \right)^2 \right\} = n \bar{\mathbf{x}}'(I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \bar{\mathbf{x}},$$

where $\bar{\mathbf{x}} = (\bar{x}_1, \dots, \bar{x}_k)'$. If we define \mathbf{x} as $\mathbf{x} = (\mathbf{x}'_1, \dots, \mathbf{x}'_k)'$, then $\mathbf{x} \sim N_{kn}(\boldsymbol{\mu}, \Omega)$ with $\boldsymbol{\mu} = \mathbf{1}_k \otimes \mu \mathbf{1}_n = \mu \mathbf{1}_{kn}$ and $\Omega = I_k \otimes \sigma^2 I_n = \sigma^2 I_{kn}$, and $\bar{\mathbf{x}} = n^{-1}(I_k \otimes \mathbf{1}'_n)\mathbf{x}$, so

$$\begin{aligned}
t_1 &= n^{-1} \mathbf{x}' (I_k \otimes \mathbf{1}_n) (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') (I_k \otimes \mathbf{1}_n') \mathbf{x} \\
&= n^{-1} \mathbf{x}' \{ (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}_n' \} \mathbf{x} = \mathbf{x}' A_1 \mathbf{x},
\end{aligned}$$

where $A_1 = n^{-1} \{ (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}_n' \}$. Since $(\mathbf{1}_n \mathbf{1}_n')^2 = n \mathbf{1}_n \mathbf{1}_n'$ and $(I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k')^2 = (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k')$, we find that A_1 is idempotent and, hence, so is $(A_1/\sigma^2)\Omega$. Thus, by Theorem 11.12, $\mathbf{x}'(A_1/\sigma^2)\mathbf{x} = t_1/\sigma^2$ has a chi-squared distribution. This distribution is central because $\lambda = \frac{1}{2} \boldsymbol{\mu}' A_1 \boldsymbol{\mu} / \sigma^2 = 0$, which follows from the fact that

$$\begin{aligned}
\{ (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}_n' \} (\mathbf{1}_k \otimes \boldsymbol{\mu} \mathbf{1}_n) &= n \boldsymbol{\mu} \{ (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \mathbf{1}_k \otimes \mathbf{1}_n \} \\
&= n \boldsymbol{\mu} \{ (\mathbf{1}_k - \mathbf{1}_k) \otimes \mathbf{1}_n \} = \mathbf{0},
\end{aligned}$$

whereas its degrees of freedom are given by

$$\begin{aligned}
r_1 &= \text{tr}\{(A_1/\sigma^2)\Omega\} = \text{tr}(A_1) = n^{-1} \text{tr}\{ (I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}_n' \} \\
&= n^{-1} \text{tr}(I_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \text{tr}(\mathbf{1}_n \mathbf{1}_n') = n^{-1} (k - 1)n = k - 1.
\end{aligned}$$

Turning to t_2 , observe that it can be written as

$$\begin{aligned}
t_2 &= \sum_{i=1}^k \left\{ \sum_{j=1}^n x_{ij}^2 - n^{-1} \left(\sum_{j=1}^n x_{ij} \right)^2 \right\} = \sum_{i=1}^k \mathbf{x}_i' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \mathbf{x}_i \\
&= \mathbf{x}' \{ I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \} \mathbf{x} = \mathbf{x}' A_2 \mathbf{x},
\end{aligned}$$

where $A_2 = I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')$. Clearly, A_2 is idempotent because $(I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')$ is idempotent. Thus, $(A_2/\sigma^2)\Omega$ is idempotent, and so $\mathbf{x}'(A_2/\sigma^2)\mathbf{x} = t_2/\sigma^2$ also has a chi-squared distribution. In particular, $t_2/\sigma^2 \sim \chi_{k(n-1)}^2$ because

$$\begin{aligned}
\text{tr}\{(A_2/\sigma^2)\Omega\} &= \text{tr}(A_2) = \text{tr}\{ I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \} \\
&= \text{tr}(I_k) \text{tr}(I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') = k(n - 1)
\end{aligned}$$

and

$$\begin{aligned}
A_2 \boldsymbol{\mu} &= \{ I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \} (\mathbf{1}_k \otimes \boldsymbol{\mu} \mathbf{1}_n) \\
&= \mathbf{1}_k \otimes \boldsymbol{\mu} (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \mathbf{1}_n = \mathbf{1}_k \otimes \boldsymbol{\mu} (\mathbf{1}_n - \mathbf{1}_n) = \mathbf{0},
\end{aligned}$$

thereby guaranteeing that $\frac{1}{2} \boldsymbol{\mu}' A_2 \boldsymbol{\mu} / \sigma^2 = 0$. Finally, we establish the independence of t_1 and t_2 by using Theorem 11.14. This simply involves verifying that $(A_1/\sigma^2)\Omega(A_2/\sigma^2) = A_1 A_2 / \sigma^2 = (\mathbf{0})$, which is an immediate consequence of the fact that

$$\mathbf{1}_n \mathbf{1}_n' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') = (\mathbf{0}).$$

Example 11.5 Let us return to the general regression model

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon},$$

where \mathbf{y} and $\boldsymbol{\epsilon}$ are $N \times 1$, X is $N \times m$, and $\boldsymbol{\beta}$ is $m \times 1$. Suppose that $\boldsymbol{\beta}$ and X are partitioned as $\boldsymbol{\beta} = (\boldsymbol{\beta}'_1 \ \boldsymbol{\beta}'_2)'$ and $X = (X_1 \ X_2)$, where $\boldsymbol{\beta}_1$ is $m_1 \times 1$, $\boldsymbol{\beta}_2$ is $m_2 \times 1$, and we wish to test the hypothesis that $\boldsymbol{\beta}_2 = \mathbf{0}$. We will assume that each component of $\boldsymbol{\beta}_2$ is estimable because this test would not be meaningful otherwise. It is easily shown that this then implies that X_2 has full column rank and $\text{rank}(X_1) = r - m_2$, where $r = \text{rank}(X)$. A test of $\boldsymbol{\beta}_2 = \mathbf{0}$ can be constructed by comparing the sum of squared errors for the reduced model $\mathbf{y} = X_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon}$, which is

$$t_1 = (\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1)'(\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1) = \mathbf{y}'(I_N - X_1(X_1'X_1)^{-}X_1')\mathbf{y},$$

with the sum of squared errors for the complete model, which is given by

$$t_2 = (\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = \mathbf{y}'(I_N - X(X'X)^{-}X')\mathbf{y}.$$

Now if $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 I_N)$, then $\mathbf{y} \sim N_N(X\boldsymbol{\beta}, \sigma^2 I_N)$. Thus, by applying Theorem 11.12 and using the fact that $X(X'X)^{-}X'X_1 = X_1$, we find that $(t_1 - t_2)/\sigma^2$ is chi-squared because

$$\begin{aligned} \left\{ \frac{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'}{\sigma^2} \right\} (\sigma^2 I_N) & \left\{ \frac{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'}{\sigma^2} \right\} \\ &= \left\{ \frac{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'}{\sigma^2} \right\}. \end{aligned}$$

In particular, if $\boldsymbol{\beta}_2 = \mathbf{0}$, $(t_1 - t_2)/\sigma^2 \sim \chi_{m_2}^2$, because

$$\begin{aligned} & \text{tr}\{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'\} \\ &= \text{tr}\{X(X'X)^{-}X'\} - \text{tr}\{X_1(X_1'X_1)^{-}X_1'\} \\ &= r - (r - m_2) = m_2 \end{aligned}$$

and

$$\begin{aligned} & \boldsymbol{\beta}'_1 X_1' \left\{ \frac{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'}{\sigma^2} \right\} X_1 \boldsymbol{\beta}_1 \\ &= \frac{\boldsymbol{\beta}'_1 X_1' X_1 \boldsymbol{\beta}_1 - \boldsymbol{\beta}'_1 X_1' X_1 \boldsymbol{\beta}_1}{\sigma^2} = 0. \end{aligned}$$

By a similar application of Theorem 11.12, we observe that $t_2/\sigma^2 \sim \chi_{N-r}^2$. In addition, it follows from Theorem 11.14 that $(t_1 - t_2)/\sigma^2$ and t_2/σ^2 are independently distributed because

$$\left\{ \frac{X(X'X)^{-}X' - X_1(X_1'X_1)^{-}X_1'}{\sigma^2} \right\} (\sigma^2 I_N) \left\{ \frac{I_N - X(X'X)^{-}X'}{\sigma^2} \right\} = 0.$$

This then permits the construction of an F statistic for testing that $\beta_2 = \mathbf{0}$; that is, if $\beta_2 = \mathbf{0}$, then the statistic

$$F = \frac{(t_1 - t_2)/m_2}{t_2/(N - r)}$$

has the F distribution with m_2 and $N - r$ degrees of freedom.

The proof Theorem 11.15, which is similar to the proof of Theorem 11.14, is left to the reader as an exercise.

Theorem 11.15 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite, and suppose that A is an $m \times m$ symmetric matrix, whereas B is an $n \times m$ matrix. If $B\Omega A = (0)$, then $\mathbf{x}'A\mathbf{x}$ and $B\mathbf{x}$ are independently distributed.

Example 11.6 Suppose that we have a random sample x_1, \dots, x_n from a normal distribution with mean μ and variance σ^2 . In Example 11.3, it was shown that $(n - 1)s^2/\sigma^2 \sim \chi_{n-1}^2$, where s^2 , the sample variance, is given by

$$s^2 = \frac{1}{(n - 1)} \sum_{i=1}^n (x_i - \bar{x})^2.$$

We will now use Theorem 11.15 to show that the sample mean,

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i,$$

is independently distributed of s^2 . In Example 11.3, we saw that s^2 is a scalar multiple of the quadratic form

$$\mathbf{x}'(I_n - n^{-1}\mathbf{1}_n\mathbf{1}_n')\mathbf{x},$$

where $\mathbf{x} = (x_1, \dots, x_n)' \sim N_n(\mu\mathbf{1}_n, \sigma^2 I_n)$. On the other hand, \bar{x} can be expressed as

$$\bar{x} = n^{-1}\mathbf{1}_n'\mathbf{x}.$$

Consequently, the independence of \bar{x} and s^2 follows from the fact that

$$\begin{aligned} \mathbf{1}_n'(\sigma^2 I_n)(I_n - n^{-1}\mathbf{1}_n\mathbf{1}_n') &= \sigma^2(\mathbf{1}_n' - n^{-1}\mathbf{1}_n'\mathbf{1}_n\mathbf{1}_n') \\ &= \sigma^2(\mathbf{1}_n' - \mathbf{1}_n') = \mathbf{0}'. \end{aligned}$$

When Ω is positive semidefinite, the condition $A\Omega B = (0)$, given in Theorem 11.14, will still guarantee that the two quadratic forms $\mathbf{x}'A\mathbf{x}$ and $\mathbf{x}'B\mathbf{x}$ are independently distributed. Likewise, when Ω is positive semidefinite, the condition $B\Omega A = (0)$, given in Theorem 11.15, will still guarantee that $\mathbf{x}'A\mathbf{x}$ and $B\mathbf{x}$ are independently

distributed. However, in these situations, a weaker set of conditions will guarantee independence. These conditions are given in Theorem 11.16 and Theorem 11.17. The proofs are left as exercises.

Theorem 11.16 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive semidefinite, and suppose that A and B are $m \times m$ symmetric matrices. Then $\mathbf{x}'A\mathbf{x}$ and $\mathbf{x}'B\mathbf{x}$ are independently distributed if

- (a) $\Omega A \Omega B \Omega = (0)$,
- (b) $\Omega A \Omega B \boldsymbol{\mu} = \mathbf{0}$,
- (c) $\Omega B \Omega A \boldsymbol{\mu} = \mathbf{0}$,
- (d) $\boldsymbol{\mu}' A \Omega B \boldsymbol{\mu} = 0$.

Theorem 11.17 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive semidefinite, and suppose that A is an $m \times m$ symmetric matrix, whereas B is an $n \times m$ matrix. If $B \Omega A \Omega = (0)$ and $B \Omega A \boldsymbol{\mu} = \mathbf{0}$, then $\mathbf{x}'A\mathbf{x}$ and $B\mathbf{x}$ are independently distributed.

Our final result can be helpful in establishing that several quadratic forms in the same normal random vector are independently distributed, with each having a chi-squared distribution.

Theorem 11.18 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite. Suppose that A_i is an $m \times m$ symmetric matrix of rank r_i , for $i = 1, \dots, k$, and $A = A_1 + \dots + A_k$ is of rank r . Consider the conditions

- (a) $A_i \Omega$ is idempotent for each i ,
- (b) $A \Omega$ is idempotent,
- (c) $A_i \Omega A_j = (0)$, for all $i \neq j$,
- (d) $r = \sum_{i=1}^k r_i$.

If any two of (a), (b), and (c) hold, or if (b) and (d) hold, then

- (i) $\mathbf{x}'A_i\mathbf{x} \sim \chi_{r_i}^2(\frac{1}{2}\boldsymbol{\mu}'A_i\boldsymbol{\mu})$,
- (ii) $\mathbf{x}'A\mathbf{x} \sim \chi_r^2(\frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu})$,
- (iii) $\mathbf{x}'A_1\mathbf{x}, \dots, \mathbf{x}'A_k\mathbf{x}$ are independently distributed.

Proof. Since Ω is positive definite, a nonsingular matrix T satisfying $\Omega = TT'$ exists, and the conditions (a)–(d) can be equivalently expressed as

- (a) $T'A_iT$ is idempotent for each i ,
- (b) $T'AT$ is idempotent,

- (c) $(T'A_iT)(T'A_jT) = (0)$, for all $i \neq j$,
 (d) $\text{rank}(T'AT) = \sum_{i=1}^k \text{rank}(T'A_iT)$.

Since $T'A_1T, \dots, T'A_kT$ and $T'AT$ satisfy the conditions of Corollary 11.8.1, we are ensured that if any two of (a), (b), and (c) hold or if (b) and (d) hold, then all four of the conditions (a)–(d) hold. Now using Theorem 11.12, (a) implies (i) and (b) implies (ii), whereas Theorem 11.14, along with (c), guarantees that (iii) holds. \square

11.6 EXPECTED VALUES OF QUADRATIC FORMS

When a quadratic form satisfies the conditions given in the theorems of Section 11.4, then its moments can be obtained directly from the appropriate chi-squared distribution. In this section, we derive formulas for means, variances, and covariances of quadratic forms that will be useful when this is not the case. We will start with the most general case in which the random vector \mathbf{x} has an arbitrary distribution. The expressions we obtain involve the matrix of second moments of \mathbf{x} , $E(\mathbf{x}\mathbf{x}')$ and the matrix of fourth moments $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')$.

Theorem 11.19 Let \mathbf{x} be an $m \times 1$ random vector having finite fourth moments, so that both $E(\mathbf{x}\mathbf{x}')$ and $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')$ exist. Denote the mean vector and covariance matrix of \mathbf{x} by $\boldsymbol{\mu}$ and Ω . If A and B are $m \times m$ symmetric matrices, then

- (a) $E(\mathbf{x}'A\mathbf{x}) = \text{tr}\{AE(\mathbf{x}\mathbf{x}')\} = \text{tr}(A\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}$,
 (b) $\text{var}(\mathbf{x}'A\mathbf{x}) = \text{tr}\{(A \otimes A)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\} - \{\text{tr}(A\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}\}^2$,
 (c) $\text{cov}(\mathbf{x}'A\mathbf{x}, \mathbf{x}'B\mathbf{x}) = \text{tr}\{(A \otimes B)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\} - \{\text{tr}(A\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}\} \times \{\text{tr}(B\Omega) + \boldsymbol{\mu}'B\boldsymbol{\mu}\}$.

Proof. The covariance matrix Ω is defined by

$$\Omega = E\{(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})'\} = E(\mathbf{x}\mathbf{x}') - \boldsymbol{\mu}\boldsymbol{\mu}',$$

so that $E(\mathbf{x}\mathbf{x}') = \Omega + \boldsymbol{\mu}\boldsymbol{\mu}'$. Since $\mathbf{x}'A\mathbf{x}$ is a scalar, we have

$$\begin{aligned} E(\mathbf{x}'A\mathbf{x}) &= E\{\text{tr}(\mathbf{x}'A\mathbf{x})\} = E\{\text{tr}(A\mathbf{x}\mathbf{x}')\} = \text{tr}\{AE(\mathbf{x}\mathbf{x}')\} \\ &= \text{tr}\{A(\Omega + \boldsymbol{\mu}\boldsymbol{\mu}')\} = \text{tr}(A\Omega) + \text{tr}(A\boldsymbol{\mu}\boldsymbol{\mu}') \\ &= \text{tr}(A\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}, \end{aligned}$$

and so (a) holds. Part (b) will follow from (c) by taking $B = A$. To prove (c), note that

$$\begin{aligned}
E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) &= E[\text{tr}\{(\mathbf{x}' \otimes \mathbf{x}')(A \otimes B)(\mathbf{x} \otimes \mathbf{x})\}] \\
&= E[\text{tr}\{(A \otimes B)(\mathbf{x} \otimes \mathbf{x})(\mathbf{x}' \otimes \mathbf{x}')\}] \\
&= \text{tr}\{(A \otimes B)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\}.
\end{aligned}$$

Applying this result, along with part (a), in the equation

$$\text{cov}(\mathbf{x}'A\mathbf{x}, \mathbf{x}'B\mathbf{x}) = E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) - E(\mathbf{x}'A\mathbf{x})E(\mathbf{x}'B\mathbf{x})$$

completes the proof. \square

When \mathbf{x} has a multivariate normal distribution, the expressions for variances and covariances, as well as for higher moments, simplify somewhat. This is a consequence of the special structure of the moments of the multivariate normal distribution. The commutation matrix K_{mm} , discussed in Chapter 8, plays a crucial role in obtaining some of these matrix expressions. We will also use the $m \times m$ matrix T_{ij} defined by

$$T_{ij} = E_{ij} + E_{ji} = \mathbf{e}_i \mathbf{e}_j' + \mathbf{e}_j \mathbf{e}_i';$$

that is, all elements of T_{ij} are equal to 0 except for the (i, j) th and (j, i) th elements, which equal 1, unless $i = j$, in which case, the only nonzero element is a 2 in the (i, i) th position. Before obtaining expressions for the variance and covariance of quadratic forms in normal variates, we will need Theorem 11.20.

Theorem 11.20 If $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$ and \mathbf{c} is a vector of constants, then

- (a) $E(\mathbf{z} \otimes \mathbf{z}) = \text{vec}(I_m)$,
- (b) $E(\mathbf{c}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}') = (0)$, $E(\mathbf{z}\mathbf{c}' \otimes \mathbf{z}\mathbf{z}') = (0)$, $E(\mathbf{z}\mathbf{z}' \otimes \mathbf{c}\mathbf{z}') = (0)$, and $E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{c}') = (0)$,
- (c) $E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}') = 2N_m + \text{vec}(I_m)\{\text{vec}(I_m)\}'$,
- (d) $\text{var}(\mathbf{z} \otimes \mathbf{z}) = 2N_m$.

Proof. Since $E(\mathbf{z}) = \mathbf{0}$, $I_m = \text{var}(\mathbf{z}) = E(\mathbf{z}\mathbf{z}')$, and so

$$E(\mathbf{z} \otimes \mathbf{z}) = E\{\text{vec}(\mathbf{z}\mathbf{z}')\} = \text{vec}\{E(\mathbf{z}\mathbf{z}')\} = \text{vec}(I_m).$$

It is easily verified using the standard normal moment generating function that

$$E(z_i^3) = 0, \quad E(z_i^4) = 3.$$

Each element of the matrices of expected values in (b) will be of the form $c_i E(z_j z_k z_l)$. Since the components of \mathbf{z} are independent, we get

$$E(z_j z_k z_l) = E(z_j)E(z_k)E(z_l) = 0$$

when the three subscripts are distinct,

$$E(z_j z_k z_l) = E(z_j^2)E(z_l) = (1)(0) = 0$$

when $j = k \neq l$, and similarly for $j = l \neq k$ and $l = k \neq j$, and

$$E(z_j z_k z_l) = E(z_j^3) = 0$$

when $j = k = l$. This proves (b). Next, we consider terms of the form $E(z_i z_j z_k z_l)$. These terms equal 1 if $i = j \neq l = k$, $i = k \neq j = l$, or $i = l \neq j = k$, equal 3 if $i = j = k = l$, and equal zero otherwise. This leads to

$$E(z_i z_j z z') = T_{ij} + \delta_{ij} I_m,$$

where δ_{ij} is the (i, j) th element of I_m . Thus,

$$\begin{aligned} E(z z' \otimes z z') &= E\left\{\left(\sum_{i=1}^m \sum_{j=1}^m E_{ij} z_i z_j\right) \otimes z z'\right\} \\ &= \sum_{i=1}^m \sum_{j=1}^m \{E_{ij} \otimes E(z_i z_j z z')\} \\ &= \sum_{i=1}^m \sum_{j=1}^m \{E_{ij} \otimes (T_{ij} + \delta_{ij} I_m)\} \\ &= \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes T_{ij}) + \sum_{i=1}^m \sum_{j=1}^m (\delta_{ij} E_{ij} \otimes I_m). \end{aligned}$$

The third result now follows because

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes T_{ij}) &= \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes E_{ji}) + \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes E_{ij}) \\ &= K_{mm} + \left\{ \sum_{i=1}^m (e_i \otimes e_i) \right\} \left\{ \sum_{j=1}^m (e'_j \otimes e'_j) \right\} \\ &= K_{mm} + \left\{ \sum_{i=1}^m \text{vec}(e_i e'_i) \right\} \left\{ \sum_{j=1}^m \{\text{vec}(e_j e'_j)\}' \right\} \\ &= K_{mm} + \text{vec}(I_m) \{\text{vec}(I_m)\}', \end{aligned}$$

$$\sum_{i=1}^m \sum_{j=1}^m (\delta_{ij} E_{ij} \otimes I_m) = \left(\sum_{i=1}^m E_{ii} \right) \otimes I_m = I_m \otimes I_m = I_{m^2},$$

and $I_{m^2} + K_{mmm} = 2N_m$. Finally, (d) is an immediate consequence of (a) and (c). \square

Theorem 11.21 generalizes the results of Theorem 11.20 to a multivariate normal distribution having a general positive definite covariance matrix.

Theorem 11.21 Let $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is positive definite, and let \mathbf{c} be an $m \times 1$ vector of constants. Then

- (a) $E(\mathbf{x} \otimes \mathbf{x}) = \text{vec}(\Omega)$,
- (b) $E(\mathbf{c}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') = (0)$, $E(\mathbf{x}\mathbf{c}' \otimes \mathbf{x}\mathbf{x}') = (0)$, $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{c}\mathbf{x}') = (0)$, and $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{c}') = (0)$,
- (c) $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') = 2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)\}'$,
- (d) $\text{var}(\mathbf{x} \otimes \mathbf{x}) = 2N_m(\Omega \otimes \Omega)$.

Proof. Let T be any nonsingular matrix satisfying $\Omega = TT'$, so that $\mathbf{z} = T^{-1}\mathbf{x}$ and $\mathbf{x} = T\mathbf{z}$, where $\mathbf{z} \sim N_m(\mathbf{0}, I_m)$. Then the results in Theorem 11.21 are consequences of Theorem 11.20 because

$$\begin{aligned} E(\mathbf{x} \otimes \mathbf{x}) &= (T \otimes T)E(\mathbf{z} \otimes \mathbf{z}) = (T \otimes T)\text{vec}(I_m) \\ &= \text{vec}(TT') = \text{vec}(\Omega), \\ E(\mathbf{c}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') &= (I_m \otimes T)E(\mathbf{c}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}')(T' \otimes T') \\ &= (I_m \otimes T)(0)(T' \otimes T') = (0) \end{aligned}$$

and

$$\begin{aligned} E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') &= (T \otimes T)E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}')(T' \otimes T') \\ &= (T \otimes T)(2N_m + \text{vec}(I_m)\{\text{vec}(I_m)\}')(T' \otimes T') \\ &= 2(T \otimes T)N_m(T' \otimes T') \\ &\quad + (T \otimes T)\text{vec}(I_m)\{\text{vec}(I_m)\}'(T' \otimes T') \\ &= 2N_m(T \otimes T)(T' \otimes T') + \text{vec}(TT')\{\text{vec}(TT')\}' \\ &= 2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)\}'. \end{aligned}$$

\square

We are now ready to obtain simplified expressions for the variance and covariance of quadratic forms in normal variates.

Theorem 11.22 Let A and B be $m \times m$ symmetric matrices, and suppose that $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is positive definite. Then

- (a) $E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) = \text{tr}(A\Omega)\text{tr}(B\Omega) + 2\text{tr}(A\Omega B\Omega),$
- (b) $\text{cov}(\mathbf{x}'A\mathbf{x}, \mathbf{x}'B\mathbf{x}) = 2\text{tr}(A\Omega B\Omega),$
- (c) $\text{var}(\mathbf{x}'A\mathbf{x}) = 2\text{tr}\{(A\Omega)^2\}.$

Proof. Since (c) is the special case of (b) in which $B = A$, we only need to prove (a) and (b). Note that by making use of Theorem 11.21, we find that

$$\begin{aligned}
 E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) &= E\{(\mathbf{x}' \otimes \mathbf{x}')(A \otimes B)(\mathbf{x} \otimes \mathbf{x})\} \\
 &= E[\text{tr}\{(A \otimes B)(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\}] \\
 &= \text{tr}\{(A \otimes B)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\} \\
 &= \text{tr}\{(A \otimes B)(2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)'\})\} \\
 &= \text{tr}\{(A \otimes B)((I_{m^2} + K_{mm})(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)'\})\} \\
 &= \text{tr}\{(A \otimes B)(\Omega \otimes \Omega)\} + \text{tr}\{(A \otimes B)K_{mm}(\Omega \otimes \Omega)\} \\
 &\quad + \text{tr}((A \otimes B)\text{vec}(\Omega)\{\text{vec}(\Omega)'\}).
 \end{aligned}$$

Now

$$\text{tr}\{(A \otimes B)(\Omega \otimes \Omega)\} = \text{tr}(A\Omega \otimes B\Omega) = \text{tr}(A\Omega)\text{tr}(B\Omega)$$

follows directly from Theorem 8.3, whereas

$$\text{tr}\{(A \otimes B)K_{mm}(\Omega \otimes \Omega)\} = \text{tr}\{(A\Omega \otimes B\Omega)K_{mm}\} = \text{tr}(A\Omega B\Omega)$$

follows from Theorem 8.26. Using the symmetry of A and Ω along with Theorem 8.10 and Theorem 8.11, the last term in $E\{\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\}$ simplifies as

$$\begin{aligned}
 \text{tr}((A \otimes B)\text{vec}(\Omega)\{\text{vec}(\Omega)'\}) &= \{\text{vec}(\Omega)\}'(A \otimes B)\text{vec}(\Omega) \\
 &= \{\text{vec}(\Omega)\}'\text{vec}(B\Omega A) = \text{tr}(A\Omega B\Omega).
 \end{aligned}$$

This then proves (a). Using the definition of covariance and Theorem 11.19(a), we also get

$$\begin{aligned}
 \text{cov}(\mathbf{x}'A\mathbf{x}, \mathbf{x}'B\mathbf{x}) &= E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) - E(\mathbf{x}'A\mathbf{x})E(\mathbf{x}'B\mathbf{x}) \\
 &= 2\text{tr}(A\Omega B\Omega),
 \end{aligned}$$

which proves (b). □

The formulas given in Theorem 11.22 become somewhat more complicated when the normal distribution has a nonnull mean vector. These formulas are given in Theorem 11.23.

Theorem 11.23 Let A and B be symmetric $m \times m$ matrices, and suppose that $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite. Then

- (a) $E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) = \text{tr}(A\Omega)\text{tr}(B\Omega) + 2\text{tr}(A\Omega B\Omega) + \text{tr}(A\Omega)\boldsymbol{\mu}'B\boldsymbol{\mu} + 4\boldsymbol{\mu}'A\Omega B\boldsymbol{\mu} + \boldsymbol{\mu}'A\boldsymbol{\mu}\text{tr}(B\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}\boldsymbol{\mu}'B\boldsymbol{\mu},$
 (b) $\text{cov}(\mathbf{x}'A\mathbf{x}, \mathbf{x}'B\mathbf{x}) = 2\text{tr}(A\Omega B\Omega) + 4\boldsymbol{\mu}'A\Omega B\boldsymbol{\mu},$
 (c) $\text{var}(\mathbf{x}'A\mathbf{x}) = 2\text{tr}\{(A\Omega)^2\} + 4\boldsymbol{\mu}'A\Omega A\boldsymbol{\mu}.$

Proof. Again, (c) is a special case of (b), so we only need to prove (a) and (b). We can write $\mathbf{x} = \mathbf{y} + \boldsymbol{\mu}$, where $\mathbf{y} \sim N_m(\mathbf{0}, \Omega)$ and, consequently,

$$\begin{aligned} E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) &= E\{(\mathbf{y} + \boldsymbol{\mu})'A(\mathbf{y} + \boldsymbol{\mu})(\mathbf{y} + \boldsymbol{\mu})'B(\mathbf{y} + \boldsymbol{\mu})\} \\ &= E\{(\mathbf{y}'A\mathbf{y} + 2\boldsymbol{\mu}'A\mathbf{y} + \boldsymbol{\mu}'A\boldsymbol{\mu})(\mathbf{y}'B\mathbf{y} + 2\boldsymbol{\mu}'B\mathbf{y} + \boldsymbol{\mu}'B\boldsymbol{\mu})\} \\ &= E(\mathbf{y}'A\mathbf{y}\mathbf{y}'B\mathbf{y}) + 2E(\mathbf{y}'A\mathbf{y}\boldsymbol{\mu}'B\mathbf{y}) + E(\mathbf{y}'A\mathbf{y})\boldsymbol{\mu}'B\boldsymbol{\mu} \\ &\quad + 2E(\boldsymbol{\mu}'A\mathbf{y}\mathbf{y}'B\mathbf{y}) + 4E(\boldsymbol{\mu}'A\mathbf{y}\boldsymbol{\mu}'B\mathbf{y}) \\ &\quad + 2E(\boldsymbol{\mu}'A\mathbf{y})\boldsymbol{\mu}'B\boldsymbol{\mu} + \boldsymbol{\mu}'A\boldsymbol{\mu}E(\mathbf{y}'B\mathbf{y}) \\ &\quad + 2\boldsymbol{\mu}'A\boldsymbol{\mu}E(\boldsymbol{\mu}'B\mathbf{y}) + \boldsymbol{\mu}'A\boldsymbol{\mu}\boldsymbol{\mu}'B\boldsymbol{\mu}. \end{aligned}$$

The sixth and eighth terms in this last expression are zero because $E(\mathbf{y}) = \mathbf{0}$, whereas it follows from Theorem 11.21(b) that the second and fourth terms are zero. To simplify the fifth term, note that

$$\begin{aligned} E(\boldsymbol{\mu}'A\mathbf{y}\boldsymbol{\mu}'B\mathbf{y}) &= E\{(\boldsymbol{\mu}'A \otimes \boldsymbol{\mu}'B)(\mathbf{y} \otimes \mathbf{y})\} = (A\boldsymbol{\mu} \otimes B\boldsymbol{\mu})'E\{(\mathbf{y} \otimes \mathbf{y})\} \\ &= \{\text{vec}(B\boldsymbol{\mu}\boldsymbol{\mu}'A)\}'\text{vec}(\Omega) = \text{tr}\{(B\boldsymbol{\mu}\boldsymbol{\mu}'A)'\Omega\} \\ &= \text{tr}(A\boldsymbol{\mu}\boldsymbol{\mu}'B\Omega) = \boldsymbol{\mu}'A\Omega B\boldsymbol{\mu}. \end{aligned}$$

Thus, using this result, Theorem 11.19(a), and Theorem 11.22(a), we find that

$$\begin{aligned} E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}) &= \text{tr}(A\Omega)\text{tr}(B\Omega) + 2\text{tr}(A\Omega B\Omega) + \text{tr}(A\Omega)\boldsymbol{\mu}'B\boldsymbol{\mu} \\ &\quad + 4\boldsymbol{\mu}'A\Omega B\boldsymbol{\mu} + \boldsymbol{\mu}'A\boldsymbol{\mu}\text{tr}(B\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu}\boldsymbol{\mu}'B\boldsymbol{\mu}, \end{aligned}$$

thereby proving (a); (b) then follows immediately from the definition of covariance and Theorem 11.19(a). \square

Example 11.7 Let us return to the subject of Example 11.4, where we defined

$$A_1 = n^{-1}\{(I_k - k^{-1}\mathbf{1}_k\mathbf{1}_k') \otimes \mathbf{1}_n\mathbf{1}_n'\}$$

and

$$A_2 = I_k \otimes (I_n - n^{-1}\mathbf{1}_n\mathbf{1}_n').$$

It was shown that if $\mathbf{x} = (\mathbf{x}'_1, \dots, \mathbf{x}'_k)'\sim N_{kn}(\boldsymbol{\mu}, \Omega)$ with $\boldsymbol{\mu} = \mathbf{1}_k \otimes \boldsymbol{\mu}_n$ and $\Omega = I_k \otimes \sigma^2 I_n$, then $t_1/\sigma^2 = \mathbf{x}'(A_1/\sigma^2)\mathbf{x} \sim \chi^2_{k-1}$ and $t_2/\sigma^2 = \mathbf{x}'(A_2/\sigma^2)\mathbf{x} \sim \chi^2_{k(n-1)}$,

independently. Since the mean of a chi-squared random variable equals its degrees of freedom, whereas the variance is two times the degrees of freedom, we can easily calculate the mean and variance of t_1 and t_2 without using the results of this section; in particular, we have

$$\begin{aligned} E(t_1) &= \sigma^2(k-1), & \text{var}(t_1) &= 2\sigma^4(k-1), \\ E(t_2) &= \sigma^2k(n-1), & \text{var}(t_2) &= 2\sigma^4k(n-1). \end{aligned}$$

Suppose now that $\mathbf{x}_i \sim N_n(\mu \mathbf{1}_n, \sigma_i^2 I_n)$, so that $\Omega = \text{var}(\mathbf{x}) = D \otimes I_n$, where $D = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$. It can be easily verified that, in this case, t_1/σ^2 and t_2/σ^2 no longer satisfy the conditions of Theorem 11.12 for chi-squaredness, but they are still independently distributed. The mean and variance of t_1 and t_2 can be computed by using Theorem 11.19 and Theorem 11.23. For instance, the mean of t_2 is given by

$$\begin{aligned} E(t_2) &= E(\mathbf{x}' A_2 \mathbf{x}) = \text{tr}(A_2 \Omega) + \boldsymbol{\mu}' A_2 \boldsymbol{\mu} \\ &= \text{tr}(\{I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')\} (D \otimes I_n)) \\ &\quad + \mu^2 (\mathbf{1}_k' \otimes \mathbf{1}_n') \{I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')\} (\mathbf{1}_k \otimes \mathbf{1}_n) \\ &= \text{tr}(D) \text{tr}(I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') + \mu^2 (\mathbf{1}_k' \mathbf{1}_k) \{\mathbf{1}_n' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \mathbf{1}_n\} \\ &= (n-1) \sum_{i=1}^k \sigma_i^2, \end{aligned}$$

whereas its variance is

$$\begin{aligned} \text{var}(t_2) &= \text{var}(\mathbf{x}' A_2 \mathbf{x}) = 2\text{tr}\{(A_2 \Omega)^2\} + 4\boldsymbol{\mu}' A_2 \Omega A_2 \boldsymbol{\mu} \\ &= 2\text{tr}\{D^2 \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')\} \\ &\quad + 4\mu^2 (\mathbf{1}_k' \otimes \mathbf{1}_n') \{D \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')\} (\mathbf{1}_k \otimes \mathbf{1}_n) \\ &= 2\text{tr}(D^2) \text{tr}\{(I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')\} \\ &\quad + 4\mu^2 (\mathbf{1}_k' D \mathbf{1}_k) \{\mathbf{1}_n' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \mathbf{1}_n\} \\ &= 2(n-1) \sum_{i=1}^k \sigma_i^4. \end{aligned}$$

We will leave it to the reader to verify that

$$\begin{aligned} E(t_1) &= (1 - k^{-1}) \sum_{i=1}^k \sigma_i^2, \\ \text{var}(t_1) &= 2 \left\{ (1 - 2k^{-1}) \sum_{i=1}^k \sigma_i^4 + k^{-2} \left(\sum_{i=1}^k \sigma_i^2 \right)^2 \right\}. \end{aligned}$$

So far we have considered the expectation of a quadratic form as well as the expectation of a product of two quadratic forms. A more general situation is one in which we need the expected value of the product of n quadratic forms. This expectation becomes more tedious to compute as n increases. For example, if A , B , and C are $m \times m$ symmetric matrices and $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, the expected value $E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x})$ can be obtained by first computing $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')$ and then applying this result in the identity

$$E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}) = \text{tr}\{(A \otimes B \otimes C)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\}.$$

The details of this derivation are left as an exercise. Magnus (1978) used an alternative method, using the cumulants of a distribution and their relationship to the moments of a distribution, to obtain the expectation of the product of an arbitrary number of quadratic forms. The results for a product of three and four quadratic forms are summarized in Theorem 11.24.

Theorem 11.24 Let A , B , C , and D be symmetric $m \times m$ matrices, and suppose that $\mathbf{x} \sim N_m(\mathbf{0}, I_m)$. Then

$$\begin{aligned} \text{(a) } E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}) &= \text{tr}(A)\text{tr}(B)\text{tr}(C) + 2\{\text{tr}(A)\text{tr}(BC) \\ &\quad + \text{tr}(B)\text{tr}(AC) + \text{tr}(C)\text{tr}(AB)\} \\ &\quad + 8\text{tr}(ABC), \end{aligned}$$

$$\begin{aligned} \text{(b) } E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}\mathbf{x}'D\mathbf{x}) &= \text{tr}(A)\text{tr}(B)\text{tr}(C)\text{tr}(D) + 8\{\text{tr}(A)\text{tr}(BCD) \\ &\quad + \text{tr}(B)\text{tr}(ACD) + \text{tr}(C)\text{tr}(ABD) \\ &\quad + \text{tr}(D)\text{tr}(ABC)\} + 4\{\text{tr}(AB)\text{tr}(CD) \\ &\quad + \text{tr}(AC)\text{tr}(BD) + \text{tr}(AD)\text{tr}(BC)\} + 2\{\text{tr}(A)\text{tr}(B)\text{tr}(CD) \\ &\quad + \text{tr}(A)\text{tr}(C)\text{tr}(BD) + \text{tr}(A)\text{tr}(D)\text{tr}(BC) \\ &\quad + \text{tr}(B)\text{tr}(C)\text{tr}(AD) + \text{tr}(B)\text{tr}(D)\text{tr}(AC) \\ &\quad + \text{tr}(C)\text{tr}(D)\text{tr}(AB)\} + 16\{\text{tr}(ABCD) \\ &\quad + \text{tr}(ABDC) + \text{tr}(ACBD)\}. \end{aligned}$$

If $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is positive definite, then A , B , C , and D appearing in the right-hand side of the equations in Theorem 11.24 are replaced by $A\Omega$, $B\Omega$, $C\Omega$, and $D\Omega$.

An alternative approach to the calculation of moments of quadratic forms uses tensor methods. This approach may be particularly appealing in those situations in

which higher ordered moments are needed or the random vector \mathbf{x} does not have a multivariate normal distribution. A detailed discussion of these tensor methods can be found in McCullagh (1987).

11.7 THE WISHART DISTRIBUTION

When x_1, \dots, x_n are independently distributed, with $x_i \sim N(0, \sigma^2)$ for every i , then

$$\mathbf{x}'\mathbf{x} = \sum_{i=1}^n x_i^2 \sim \sigma^2 \chi_n^2,$$

where $\mathbf{x}' = (x_1, \dots, x_n)$; that is, $\mathbf{x}'\mathbf{x}/\sigma^2$ has a chi-squared distribution with n degrees of freedom. A natural matrix generalization of this situation, one that has important applications in multivariate analysis, involves the distribution of

$$X'X = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i',$$

where $X' = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ is an $m \times n$ matrix, such that $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent and $\mathbf{x}_i \sim N_m(\mathbf{0}, \Omega)$ for each i . Thus, the components of the j th column of X are independently distributed each as $N(0, \sigma_{jj})$, where σ_{jj} is the j th diagonal element of Ω , so that the j th diagonal element of $X'X$ has the distribution $\sigma_{jj}^2 \chi_n^2$. The joint distribution of all elements of the $m \times m$ matrix $X'X$ is called the Wishart distribution with scale matrix Ω and degrees of freedom n , and it will be denoted by $W_m(\Omega, n)$. This Wishart distribution, like the chi-squared distribution χ_n^2 , is said to be central. More generally, if $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent and $\mathbf{x}_i \sim N_m(\boldsymbol{\mu}_i, \Omega)$, then $X'X$ has the noncentral Wishart distribution with noncentrality matrix $\Phi = \frac{1}{2}M'M$, where M' is the $m \times n$ matrix given by $M' = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n)$. We will denote this noncentral Wishart distribution as $W_m(\Omega, n, \Phi)$. Additional information regarding the Wishart distribution, such as the form of its density function, can be found in texts on multivariate analysis such as Srivastava and Khatri (1979) and Muirhead (1982).

If A is an $n \times n$ symmetric matrix and X' is an $m \times n$ matrix, then the matrix $X'AX$ is sometimes called a generalized quadratic form. Theorem 11.25 gives some generalizations of the results obtained in Section 11.4 and Section 11.5 regarding quadratic forms to these generalized quadratic forms.

Theorem 11.25 Let X' be an $m \times n$ matrix whose columns are independently distributed, with the i th column having the $N_m(\boldsymbol{\mu}_i, \Omega)$ distribution, where Ω is positive definite. Suppose that A and B are $n \times n$ symmetric matrices whereas C is $k \times n$. Let $M' = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n)$, $\Phi = \frac{1}{2}M'AM$, and $r = \text{rank}(A)$. Then

- (a) $X'AX \sim W_m(\Omega, r, \Phi)$, if A is idempotent,
- (b) $X'AX$ and $X'BX$ are independently distributed if $AB = (0)$,
- (c) $X'AX$ and CX are independently distributed if $CA = (0)$.

Proof. The proof of (a) will be complete if we can show that an $m \times r$ matrix Y' exists, such that $X'AX = Y'Y$, where the columns of Y' are independently distributed, each having a normal distribution with the same covariance matrix Ω , and $\frac{1}{2}E(Y')E(Y) = \Phi$. Since the columns of X' are independently distributed, it follows that

$$\text{vec}(X') \sim N_{nm}(\text{vec}(M'), I_n \otimes \Omega).$$

Since A is symmetric, idempotent, and has rank r , an $n \times r$ matrix P must exist which satisfies $A = PP'$ and $P'P = I_r$. Consequently, $X'AX = Y'Y$, where the $m \times r$ matrix $Y' = X'P$, so that

$$\begin{aligned} \text{vec}(Y') &= \text{vec}(X'P) = (P' \otimes I_m)\text{vec}(X') \\ &\sim N_{mr}((P' \otimes I_m)\text{vec}(M'), (P' \otimes I_m)(I_n \otimes \Omega)(P \otimes I_m)) \\ &\sim N_{mr}(\text{vec}(M'P), (I_r \otimes \Omega)), \end{aligned}$$

which means that the columns of Y' are independently and normally distributed, each with covariance matrix Ω . Furthermore,

$$\frac{1}{2}E(Y')E(Y) = \frac{1}{2}M'PP'M = \frac{1}{2}M'AM = \Phi,$$

and so (a) follows. To prove (b), note that because A and B are symmetric, $AB = (0)$ implies that $AB = BA$, so A and B are diagonalized by the same orthogonal matrix; that is, there exist diagonal matrices C and D and an orthogonal matrix Q , such that $Q'AQ = C$ and $Q'BQ = D$. Furthermore, $AB = (0)$ implies that $CD = (0)$, so that by appropriately choosing Q , we will have $C = \text{diag}(c_1, \dots, c_h, 0, \dots, 0)$ and $D = \text{diag}(0, \dots, 0, d_{h+1}, \dots, d_n)$ for some h . Thus, if we let $U = Q'X$, we find that

$$X'AX = U'CU = \sum_{i=1}^h c_i \mathbf{u}_i \mathbf{u}_i', \quad X'BX = U'DU = \sum_{i=h+1}^n d_i \mathbf{u}_i \mathbf{u}_i',$$

where \mathbf{u}_i is the i th column of U' . Since $\text{vec}(U') \sim N_{nm}(\text{vec}(M'Q), (I_n \otimes \Omega))$, these columns are independently distributed, and so (b) follows. The proof of (c) is similar to that of (b). \square

An application of our next result indicates that a principal submatrix of a Wishart matrix also has a Wishart distribution.

Theorem 11.26 Suppose that $V \sim W_m(\Omega, n, \Phi)$ and A is a $p \times m$ matrix of constants with $\text{rank}(A) = p$. Then $AV A' \sim W_p(A\Omega A', n, A\Phi A')$.

Proof. Since $V \sim W_m(\Omega, n, \Phi)$, it can be written as $V = X'X$, where the columns of X' are independently distributed with the i th column $\mathbf{x}_i \sim N_m(\boldsymbol{\mu}_i, \Omega)$

and the matrix M' having μ_i as its i th column satisfies $\frac{1}{2}M'M = \Phi$. Let $Y' = AX'$, so that the columns of Y' are also independent with i th column $y_i = Ax_i \sim N_p(A\mu_i, A\Omega A')$. It follows from Theorem 11.25 that $AV A' = AX'XA' = Y'Y$ has the distribution $W_p(A\Omega A', n, \Phi_*)$. The matrix Φ_* satisfies

$$\Phi_* = \frac{1}{2}E(Y')E(Y) = \frac{1}{2}AE(X')E(X)A' = \frac{1}{2}AM'MA' = A\Phi A',$$

and so the proof is complete. \square

If a matrix V having a Wishart distribution is partitioned in the form

$$V = \begin{bmatrix} V_{11} & V_{12} \\ V_{12}' & V_{22} \end{bmatrix},$$

where V_{11} and V_{22} are square matrices, then it is an immediate consequence of Theorem 11.26 that V_{11} , as well as V_{22} , has a Wishart distribution. Theorem 11.27 indicates that the Schur complement of V_{11} in V also has a Wishart distribution.

Theorem 11.27 Suppose that $V \sim W_m(\Omega, n)$, where Ω is positive definite. Partition V and Ω as

$$V = \begin{bmatrix} V_{11} & V_{12} \\ V_{12}' & V_{22} \end{bmatrix}, \quad \Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega_{12}' & \Omega_{22} \end{bmatrix},$$

where V_{11} and Ω_{11} are $m_1 \times m_1$ and V_{22} and Ω_{22} are $m_2 \times m_2$. Then $V_{11} - V_{12}V_{22}^{-1}V_{12}' \sim W_{m_1}(\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{12}', n - m_2)$.

Proof. Since $V \sim W_m(\Omega, n)$, it can be expressed as $V = X'X$, where the columns of X' are independently distributed each having the distribution $N_m(\mathbf{0}, \Omega)$. Partitioning the $n \times m$ matrix X as $X = (X_1, X_2)$, where X_1 is $n \times m_1$, we find that $V_{11} = X_1'X_1$, $V_{22} = X_2'X_2$, and $V_{12} = X_1'X_2$. Thus,

$$\begin{aligned} V_{11} - V_{12}V_{22}^{-1}V_{12}' &= X_1'X_1 - X_1'X_2(X_2'X_2)^{-1}X_2'X_1 \\ &= X_1'\{I_n - X_2(X_2'X_2)^{-1}X_2'\}X_1 = X_1'AX_1, \end{aligned}$$

where $A = I_n - X_2(X_2'X_2)^{-1}X_2$. Now, from Example 7.3, we know that, given X_2 , the columns of X_1' are independently and normally distributed with covariance matrix $\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{12}'$, whereas $E(X_1'|X_2) = \Omega_{12}\Omega_{22}^{-1}X_2'$. Since A is a symmetric idempotent matrix of rank $n - m_2$, it follows from Theorem 11.25 that, given X_2 , $X_1'AX_1 = V_{11} - V_{12}V_{22}^{-1}V_{12}' \sim W_{m_1}(\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{12}', n - m_2)$. This Wishart distribution is central because

$$\begin{aligned} E(X_1'|X_2)AE(X_1|X_2) &= \Omega_{12}\Omega_{22}^{-1}X_2'\{I_n - X_2(X_2'X_2)^{-1}X_2'\}X_2\Omega_{22}^{-1}\Omega_{12}' \\ &= \Omega_{12}\Omega_{22}^{-1}\{X_2'X_2 - X_2'X_2\}\Omega_{22}^{-1}\Omega_{12}' \\ &= (0). \end{aligned}$$

The result now follows because this conditional distribution of $X_1'AX_1$ does not depend on X_2 . \square

If the columns of the $m \times n$ matrix X' are independent and identically distributed as $N_m(\mathbf{0}, \Omega)$ and M' is an $m \times n$ matrix of constants, then $V = (X + M)'(X + M)$ has the Wishart distribution $W_m(\Omega, n, \frac{1}{2}M'M)$. A more general situation is one in which the columns of X' are independent and identically distributed having zero mean vector and some nonnormal multivariate distribution. In this case, the distribution of $V = (X + M)'(X + M)$, which may be complicated, will depend on the specific nonnormal distribution. In particular, the moments of V are directly related to the moments of the columns of X' . Our next result gives expressions for the first two moments of V when $M = (\mathbf{0})$. Since V is a matrix and joint distributions are more conveniently handled in the form of vectors, we will vectorize V ; that is, for instance, variances and covariances of the elements of V can be obtained from the matrix $\text{var}\{\text{vec}(V)\}$.

Theorem 11.28 Let the columns of the $m \times n$ matrix $X' = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ be independently and identically distributed with $E(\mathbf{x}_i) = \mathbf{0}$, $\text{var}(\mathbf{x}_i) = \Omega$, and $E(\mathbf{x}_i\mathbf{x}_i' \otimes \mathbf{x}_i\mathbf{x}_i') = \Psi$. If $V = X'X$, then

- (a) $E(V) = n\Omega$,
- (b) $\text{var}\{\text{vec}(V)\} = n\{\Psi - \text{vec}(\Omega)\text{vec}(\Omega)'\}$.

Proof. Since $E(\mathbf{x}_i) = \mathbf{0}$, we have $\Omega = E(\mathbf{x}_i\mathbf{x}_i')$, and so

$$E(V) = E(X'X) = \sum_{i=1}^n E(\mathbf{x}_i\mathbf{x}_i') = \sum_{i=1}^n \Omega = n\Omega.$$

In addition, because $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent, we have

$$\begin{aligned} \text{var}\{\text{vec}(V)\} &= \text{var}\left\{\text{vec}\left(\sum_{i=1}^n \mathbf{x}_i\mathbf{x}_i'\right)\right\} = \text{var}\left\{\sum_{i=1}^n \text{vec}(\mathbf{x}_i\mathbf{x}_i')\right\} \\ &= \sum_{i=1}^n \text{var}\{\text{vec}(\mathbf{x}_i\mathbf{x}_i')\} = \sum_{i=1}^n \text{var}(\mathbf{x}_i \otimes \mathbf{x}_i) \\ &= \sum_{i=1}^n \{E(\mathbf{x}_i\mathbf{x}_i' \otimes \mathbf{x}_i\mathbf{x}_i') - E(\mathbf{x}_i \otimes \mathbf{x}_i)E(\mathbf{x}_i' \otimes \mathbf{x}_i')\} \\ &= \sum_{i=1}^n \{\Psi - \text{vec}(\Omega)\text{vec}(\Omega)'\} = n\{\Psi - \text{vec}(\Omega)\text{vec}(\Omega)'\}, \end{aligned}$$

and so the proof is complete. \square

The expression for $\text{var}\{\text{vec}(V)\}$ simplifies when V has a Wishart distribution because of the special structure of the fourth moments of the normal distribution. This simplified expression is given in Theorem 11.29. Note that although this theorem is stated for normally distributed columns, the first result given applies to the general case as well.

Theorem 11.29 Let the columns of the $m \times n$ matrix X' be independently and identically distributed as $N_m(\mathbf{0}, \Omega)$. Define $V = (X + M)'(X + M)$, where $M' = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n)$ is an $m \times n$ matrix of constants, so that $V \sim W_m(\Omega, n, \frac{1}{2}M'M)$. Then

- (a) $E(V) = n\Omega + M'M$,
- (b) $\text{var}\{\text{vec}(V)\} = 2N_m\{n(\Omega \otimes \Omega) + \Omega \otimes M'M + M'M \otimes \Omega\}$.

Proof. Since $E(X) = (0)$ and $E(X'X) = n\Omega$ from Theorem 11.28, it follows that

$$\begin{aligned} E(V) &= E(X'X + X'M + M'X + M'M) \\ &= E(X'X) + M'M = n\Omega + M'M. \end{aligned}$$

Proceeding as in the proof of Theorem 11.28, we obtain

$$\text{var}\{\text{vec}(V)\} = \sum_{i=1}^n \text{var}\{(\mathbf{x}_i + \boldsymbol{\mu}_i) \otimes (\mathbf{x}_i + \boldsymbol{\mu}_i)\}. \quad (11.14)$$

However,

$$\begin{aligned} (\mathbf{x}_i + \boldsymbol{\mu}_i) \otimes (\mathbf{x}_i + \boldsymbol{\mu}_i) &= \mathbf{x}_i \otimes \mathbf{x}_i + \mathbf{x}_i \otimes \boldsymbol{\mu}_i + \boldsymbol{\mu}_i \otimes \mathbf{x}_i + \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \\ &= \mathbf{x}_i \otimes \mathbf{x}_i + (I_{m^2} + K_{mm})(\mathbf{x}_i \otimes \boldsymbol{\mu}_i) + \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i \\ &= \mathbf{x}_i \otimes \mathbf{x}_i + 2N_m(I_m \otimes \boldsymbol{\mu}_i)\mathbf{x}_i + \boldsymbol{\mu}_i \otimes \boldsymbol{\mu}_i. \end{aligned}$$

Since all first- and third-order moments of \mathbf{x}_i are equal to 0, $\mathbf{x}_i \otimes \mathbf{x}_i$ and \mathbf{x}_i are uncorrelated, and so using Theorem 11.21 and Problem 8.60, we find that

$$\begin{aligned} \text{var}\{(\mathbf{x}_i + \boldsymbol{\mu}_i) \otimes (\mathbf{x}_i + \boldsymbol{\mu}_i)\} &= \text{var}(\mathbf{x}_i \otimes \mathbf{x}_i) + \text{var}\{2N_m(I_m \otimes \boldsymbol{\mu}_i)\mathbf{x}_i\} \\ &= 2N_m(\Omega \otimes \Omega) \\ &\quad + 4N_m(I_m \otimes \boldsymbol{\mu}_i)\Omega(I_m \otimes \boldsymbol{\mu}_i')N_m \\ &= 2N_m(\Omega \otimes \Omega) + 4N_m(\Omega \otimes \boldsymbol{\mu}_i\boldsymbol{\mu}_i')N_m \\ &= 2N_m(\Omega \otimes \Omega + \Omega \otimes \boldsymbol{\mu}_i\boldsymbol{\mu}_i' \\ &\quad + \boldsymbol{\mu}_i\boldsymbol{\mu}_i' \otimes \Omega). \end{aligned} \quad (11.15)$$

Now substituting (11.15) in (11.14) and simplifying, we obtain (b). \square

Example 11.8 In Example 11.3 and Example 11.6, it was shown that, when sampling from a normal distribution, a constant multiple of the sample variance s^2 has a chi-squared distribution, and it is independently distributed of the sample mean \bar{x} . In this example, we consider the multivariate version of this problem involving \bar{x} and S ; that is, suppose that $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independently distributed with $\mathbf{x}_i \sim N_m(\boldsymbol{\mu}, \Omega)$ for each i , and define X' to be the $m \times n$ matrix $(\mathbf{x}_1, \dots, \mathbf{x}_n)$. Then the sample mean vector and sample covariance matrix can be expressed as

$$\bar{\mathbf{x}} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i = \frac{1}{n} X' \mathbf{1}_n$$

and

$$\begin{aligned} S &= \frac{1}{n-1} \sum_{i=1}^n (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})' = \frac{1}{n-1} \left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i' - n \bar{\mathbf{x}} \bar{\mathbf{x}}' \right) \\ &= \frac{1}{n-1} (X'X - n^{-1} X' \mathbf{1}_n \mathbf{1}_n' X) = \frac{1}{n-1} X' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') X. \end{aligned}$$

Since $A = (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n')$ is idempotent and $\text{rank}(A) = \text{tr}(A) = n-1$, it follows from Theorem 11.25(a) that $(n-1)S$ has a Wishart distribution. To determine its noncentrality matrix, note that $M' = (\boldsymbol{\mu}, \dots, \boldsymbol{\mu}) = \boldsymbol{\mu} \mathbf{1}_n'$, so that

$$M'AM = \boldsymbol{\mu} \mathbf{1}_n' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') \mathbf{1}_n \boldsymbol{\mu}' = \boldsymbol{\mu} (n - n) \boldsymbol{\mu}' = (0).$$

Thus, $(n-1)S$ has the central Wishart distribution $W_m(\Omega, n-1)$. Furthermore, using Theorem 11.25(c), we see that S and $\bar{\mathbf{x}}$ are independently distributed because

$$\mathbf{1}_n' (I_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') = (\mathbf{1}_n' - \mathbf{1}_n') = \mathbf{0}'.$$

In addition, it follows from Theorem 11.28 and Theorem 11.21 that

$$E(S) = \Omega, \quad \text{var}\{\text{vec}(S)\} = \frac{2}{n-1} N_m(\Omega \otimes \Omega) = \frac{2}{n-1} N_m(\Omega \otimes \Omega) N_m.$$

The redundant elements in $\text{vec}(S)$ can be eliminated by using $\mathbf{v}(S)$. Since $\mathbf{v}(S) = D_m^+ \text{vec}(S)$, where D_m is the duplication matrix discussed in Section 8.7, we have

$$\text{var}\{\mathbf{v}(S)\} = \frac{2}{n-1} D_m^+ N_m(\Omega \otimes \Omega) N_m D_m^+.$$

In some situations, we may be interested only in the sample variances and not the sample covariances; that is, the random vector of interest here is the $m \times 1$ vector $\mathbf{s} = (s_{11}, \dots, s_{mm})'$. Expressions for the mean vector and covariance matrix of \mathbf{s}

are easily obtained from the formulas above because $\mathbf{s} = \mathbf{w}(S) = \Psi_m \text{vec}(S)$ as seen in Problem 8.48, where

$$\Psi_m = \sum_{i=1}^m \mathbf{e}_{i,m} (\mathbf{e}_{i,m} \otimes \mathbf{e}_{i,m})'.$$

Thus, using the properties of Ψ_m obtained in Problem 8.48, we find that

$$\begin{aligned} E(\mathbf{s}) &= \Psi_m \text{vec}\{E(S)\} = \Psi_m \text{vec}(\Omega) = \mathbf{w}(\Omega), \\ \text{var}(\mathbf{s}) &= \Psi_m \text{var}\{\text{vec}(S)\} \Psi_m' = \Psi_m \left\{ \frac{2}{n-1} N_m(\Omega \otimes \Omega) N_m \right\} \Psi_m' \\ &= \frac{2}{n-1} \Psi_m (\Omega \otimes \Omega) \Psi_m' = \frac{2}{n-1} (\Omega \odot \Omega), \end{aligned}$$

where \odot is the Hadamard product.

Example 11.9 We can use the perturbation formulas for eigenvalues and eigenvectors of a symmetric matrix obtained in Section 9.6 to approximate the distributions of an eigenvalue or an eigenvector of a matrix having a Wishart distribution. One important application in statistics that uses these asymptotic distributions is principal components analysis, an analysis involving the eigenvalues and eigenvectors of the $m \times m$ sample covariance matrix S . The exact distributions of an eigenvalue and an eigenvector of S are rather complicated, whereas their asymptotic distributions follow in a fairly straightforward manner from the asymptotic distribution of S . Now it can be shown by using the central limit theorem (see Muirhead, 1982) that $\sqrt{n-1} \text{vec}(S)$ has an asymptotic normal distribution. In particular, using results from Example 11.8, we have, asymptotically,

$$\sqrt{n-1} \{\text{vec}(S) - \text{vec}(\Omega)\} \sim N_{m^2}(\mathbf{0}, 2N_m(\Omega \otimes \Omega)),$$

where Ω is the population covariance matrix. Let $W = S - \Omega$ and $W_* = \sqrt{n-1}W$, so that $\text{vec}(W_*)$ has the asymptotic normal distribution indicated above. Suppose that γ_i is a normalized eigenvector of $S = \Omega + W$ corresponding to the i th largest eigenvalue λ_i , whereas \mathbf{q}_i is a normalized eigenvector of Ω corresponding to its i th largest eigenvalue x_i . Now if x_i is a distinct eigenvalue of Ω , then we have the first-order approximations from Section 9.6

$$\begin{aligned} \lambda_i &= x_i + \mathbf{q}_i' W \mathbf{q}_i = x_i + (\mathbf{q}_i' \otimes \mathbf{q}_i') \text{vec}(W), \\ \gamma_i &= \mathbf{q}_i - (\Omega - x_i I_m)^+ W \mathbf{q}_i \\ &= \mathbf{q}_i - \{\mathbf{q}_i' \otimes (\Omega - x_i I_m)^+\} \text{vec}(W). \end{aligned} \tag{11.16}$$

Thus, the asymptotic normality of $a_i = \sqrt{n-1}(\lambda_i - x_i)$ follows from the asymptotic normality of $\text{vec}(W_*)$. Furthermore, we have, asymptotically,

$$\begin{aligned} E(a_i) &= (\mathbf{q}'_i \otimes \mathbf{q}'_i) E\{\text{vec}(W_*)\} = (\mathbf{q}'_i \otimes \mathbf{q}'_i) \mathbf{0} = 0, \\ \text{var}(a_i) &= (\mathbf{q}'_i \otimes \mathbf{q}'_i) (\text{var}\{\text{vec}(W_*)\}) (\mathbf{q}_i \otimes \mathbf{q}_i) \\ &= (\mathbf{q}'_i \otimes \mathbf{q}'_i) (2N_m(\Omega \otimes \Omega)) (\mathbf{q}_i \otimes \mathbf{q}_i) \\ &= 2(\mathbf{q}'_i \Omega \mathbf{q}_i \otimes \mathbf{q}'_i \Omega \mathbf{q}_i) = 2x_i^2; \end{aligned}$$

that is, for large n , $\lambda_i \sim N(x_i, 2x_i^2/(n-1))$, approximately. Similarly, $\mathbf{b}_i = \sqrt{n-1}(\gamma_i - \mathbf{q}_i)$ is asymptotically normal with

$$\begin{aligned} E(\mathbf{b}_i) &= -\{\mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} E\{\text{vec}(W_*)\} \\ &= -\{\mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} \mathbf{0} = \mathbf{0}, \\ \Xi &= \text{var}(\mathbf{b}_i) = \{\mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} \{\text{var}\{\text{vec}(W_*)\}\} \\ &\quad \times \{\mathbf{q}_i \otimes (\Omega - x_i I_m)^+\}' \\ &= \{\mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} \{2N_m(\Omega \otimes \Omega)\} \{\mathbf{q}_i \otimes (\Omega - x_i I_m)^+\} \\ &= \{(\Omega - x_i I_m)^+ \otimes \mathbf{q}'_i + \mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} (\Omega \otimes \Omega) \\ &\quad \times \{\mathbf{q}_i \otimes (\Omega - x_i I_m)^+\} \\ &= \mathbf{q}'_i \Omega \mathbf{q}_i \otimes (\Omega - x_i I_m)^+ \Omega (\Omega - x_i I_m)^+ \\ &= x_i \left\{ \sum_{j \neq i} \frac{x_j}{(x_j - x_i)^2} \mathbf{q}_j \mathbf{q}'_j \right\}, \end{aligned}$$

and so for large n , $\gamma_i \sim N_m(\mathbf{q}_i, (n-1)^{-1}\Xi)$, approximately. While we can use the first-order approximations in (11.16) to obtain the asymptotic distributions, we can use higher order approximations, such as those given in Theorem 9.5, to further improve the performance of these asymptotic distributions. The most common application of this process involves asymptotic chi-squared distributions, so we will illustrate the basic idea with the statistic

$$t = \frac{(n-1)(\lambda_i - x_i)^2}{2x_i^2},$$

which, because of the asymptotic normality of λ_i , is asymptotically chi-squared with one degree of freedom. The mean of this chi-squared distribution is 1, whereas the exact mean of t is of the form

$$E(t) = 1 + \sum_{j=1}^{\infty} \frac{c_j}{(n-1)^{(j+1)/2}},$$

where the c_j 's are constants. We can use the higher order approximations of λ_i to determine the first constant c_1 , and then this may be used to compute an adjusted statistic

$$t_* = \left\{ 1 - \frac{c_1}{(n-1)} \right\} t.$$

The mean of this adjusted statistic is

$$\begin{aligned} E(t_*) &= \left\{ 1 - \frac{c_1}{(n-1)} \right\} E(t) \\ &= \left\{ 1 - \frac{c_1}{(n-1)} \right\} \left(1 + \sum_{j=1}^{\infty} \frac{c_j}{(n-1)^{(j+1)/2}} \right) \\ &= 1 + \sum_{j=2}^{\infty} \frac{d_j}{(n-1)^{(j+1)/2}}, \end{aligned}$$

where the d_j 's are constants that are functions of the c_j 's. Note that the mean of t_* converges to 1 at a faster rate than does $E(t)$. For this reason, the chi-squared distribution with one degree of freedom should approximate the distribution of this adjusted statistic better than it would approximate the distribution of t . This type of adjustment of asymptotically chi-squared statistics is commonly referred to as a Bartlett adjustment (Bartlett, 1937, 1947). Some further discussion of Bartlett adjustments can be found in Barndorff-Nielsen and Cox (1994).

Some of the inequalities for eigenvalues developed in Chapter 3 have important applications regarding the distributions of eigenvalues of certain functions of Wishart matrices. One such application is illustrated in Example 11.10.

Example 11.10 A multivariate analysis of variance, such as the multivariate one-way classification model discussed in Example 3.16, uses the eigenvalues of BW^{-1} , where the $m \times m$ matrices B and W are independently distributed with $B \sim W_m(I_m, b, \Phi)$ and $W \sim W_m(I_m, w)$ (Problem 11.49). We will show that if the rank of the noncentrality matrix Φ is $r < m$ and V_1 and V_2 are independently distributed with $V_1 \sim W_{m-r}(I_{m-r}, b-r)$ and $V_2 \sim W_{m-r}(I_{m-r}, w)$, then

$$P\{\lambda_{r+i}(BW^{-1}) > c\} \leq P\{\lambda_i(V_1V_2^{-1}) > c\},$$

for $i = 1, \dots, m-r$ and any constant c . This result is useful in determining the dimensionality in a canonical variate analysis (see Schott, 1984). Since $\text{rank}(\Phi) = r$, an $r \times m$ matrix T exists, such that $\frac{1}{2}T'T = \Phi$. If we define the $m \times b$ matrix $M' = (T' \quad 0)$, then because $\frac{1}{2}M'M = \Phi$ and $B \sim W_m(I_m, b, \Phi)$, it follows that B can be expressed as $B = X'X$, where X' is an $m \times b$ matrix for which

$\text{vec}(X') \sim N_{bm}(\text{vec}(M'), I_b \otimes I_m)$. Partitioning X' as $X' = (X'_1 \ X'_2)$, where X'_1 is $m \times r$, we find that

$$B = X'_1 X_1 + X'_2 X_2 = B_1 + B_2,$$

where $B_1 = X'_1 X_1 \sim W_m(I_m, r, \Phi)$ and $B_2 = X'_2 X_2 \sim W_m(I_m, b - r)$ because

$$\text{vec}(X'_1) \sim N_{rm}(\text{vec}(T'), I_r \otimes I_m)$$

and

$$\text{vec}(X'_2) \sim N_{(b-r)m}(\text{vec}\{(0)\}, I_{b-r} \otimes I_m).$$

Now for fixed B_1 , let F be any $m \times (m - r)$ matrix satisfying $F' B_1 F = (0)$ and $F' F = I_{m-r}$, and define the sets

$$\begin{aligned} S_1(B_1) &= \{B_2, W : \lambda_{r+i}(BW^{-1}) > c\}, \\ S_2(B_1) &= \{B_2, W : \lambda_i\{(F' B_2 F)(F' W F)^{-1}\} > c\}. \end{aligned}$$

It follows from Problem 3.47(a) that

$$\lambda_i\{(F' B F)(F' W F)^{-1}\} = \lambda_i\{(F' B_2 F)(F' W F)^{-1}\} \geq \lambda_{r+i}(BW^{-1}),$$

so for each fixed B_1 , $S_1(B_1) \subseteq S_2(B_1)$, and it follows from Theorem 11.26 that $V_1 = F' B_2 F \sim W_{m-r}(I_{m-r}, b - r)$ and $V_2 = F' W F \sim W_{m-r}(I_{m-r}, w)$. Consequently, if $g(W)$, $f_1(B_1)$, and $f_2(B_2)$ are the density functions for W , B_1 , and B_2 , respectively, then

$$\begin{aligned} \int_{S_1(B_1)} g(W) f_2(B_2) dW dB_2 &\leq \int_{S_2(B_1)} g(W) f_2(B_2) dW dB_2 \\ &= P\{\lambda_i(V_1 V_2^{-1}) > c\}. \end{aligned}$$

If we also define the sets

$$\begin{aligned} C_1 &= \{B_1, B_2, W : \lambda_{r+i}(BW^{-1}) > c\}, \\ C_2 &= \{B_1 : B_1 \text{ positive definite}\}, \end{aligned}$$

then the desired result follows because

$$\begin{aligned} P\{\lambda_{r+i}(BW^{-1}) > c\} &= \int_{C_1} g(W) f_1(B_1) f_2(B_2) dW dB_1 dB_2 \\ &= \int_{C_2} \left\{ \int_{S_1(B_1)} g(W) f_2(B_2) dW dB_2 \right\} f_1(B_1) dB_1 \\ &\leq \int_{C_2} P\{\lambda_i(V_1 V_2^{-1}) > c\} f_1(B_1) dB_1 \\ &= P\{\lambda_i(V_1 V_2^{-1}) > c\}. \end{aligned}$$

We can use the relationship between the sample correlation and sample covariance matrices and the expression for $\text{var}\{\text{vec}(S)\}$ given in Example 11.8 to obtain an expression for the asymptotic covariance matrix of $\text{vec}(R)$. This is the subject of our final example.

Example 11.11 As in Example 11.8, let $\mathbf{x}_1, \dots, \mathbf{x}_n$ be independently distributed with $\mathbf{x}_i \sim N_m(\boldsymbol{\mu}, \Omega)$, for each i , and let S and R be the sample covariance and correlation matrices computed from this sample. Thus, if we use the notation $D_X^a = \text{diag}(x_{11}^a, \dots, x_{mm}^a)$, where X is an $m \times m$ matrix, then the sample correlation matrix can be expressed as

$$R = D_S^{-1/2} S D_S^{-1/2},$$

whereas the population correlation matrix is given by

$$P = D_\Omega^{-1/2} \Omega D_\Omega^{-1/2}.$$

Note that if we define $\mathbf{y}_i = D_\Omega^{-1/2} \mathbf{x}_i$, then $\mathbf{y}_1, \dots, \mathbf{y}_n$ are independently distributed with $\mathbf{y}_i \sim N_m(D_\Omega^{-1/2} \boldsymbol{\mu}, P)$. If S_* is the sample covariance matrix computed from the \mathbf{y}_i 's, then $S_* = D_\Omega^{-1/2} S D_\Omega^{-1/2}$, $D_{S_*}^{-1/2} = D_S^{-1/2} D_\Omega^{1/2} = D_\Omega^{1/2} D_S^{-1/2}$, and so

$$\begin{aligned} D_{S_*}^{-1/2} S_* D_{S_*}^{-1/2} &= D_S^{-1/2} D_\Omega^{1/2} (D_\Omega^{-1/2} S D_\Omega^{-1/2}) D_\Omega^{1/2} D_S^{-1/2} \\ &= D_S^{-1/2} S D_S^{-1/2} = R; \end{aligned}$$

that is, the sample correlation matrix computed from the \mathbf{y}_i 's is the same as that computed from the \mathbf{x}_i 's. If $A = S_* - P$, then the first-order approximation for R is given by (see Problem 9.23)

$$R = P + A - \frac{1}{2}(P D_A + D_A P),$$

and so

$$\begin{aligned} \text{vec}(R) &= \text{vec}(P) + \text{vec}(A) - \frac{1}{2}\{\text{vec}(P D_A) + \text{vec}(D_A P)\} \\ &= \text{vec}(P) + \text{vec}(A) - \frac{1}{2}\{(I_m \otimes P) + (P \otimes I_m)\}\text{vec}(D_A) \\ &= \text{vec}(P) + \left(I_{m^2} - \frac{1}{2}\{(I_m \otimes P) \right. \\ &\quad \left. + (P \otimes I_m)\} \Lambda_m \right) \text{vec}(A), \end{aligned} \tag{11.17}$$

where

$$\Lambda_m = \sum_{i=1}^m (E_{ii} \otimes E_{ii}).$$

Thus, because

$$\text{var}\{\text{vec}(A)\} = \text{var}\{\text{vec}(S_*)\} = \frac{2}{n-1}N_m(P \otimes P)N_m,$$

we get the first-order approximation

$$\text{var}\{\text{vec}(R)\} = \frac{2}{n-1}HN_m(P \otimes P)N_mH',$$

where the matrix H is the premultiplier on $\text{vec}(A)$ in the last expression given in (11.17). Simplification (see Problem 11.53) leads to

$$\text{var}\{\text{vec}(R)\} = \frac{2}{n-1}N_m\Theta N_m, \quad (11.18)$$

where

$$\Theta = \{I_{m^2} - (I_m \otimes P)\Lambda_m\}(P \otimes P)\{I_{m^2} - \Lambda_m(I_m \otimes P)\}.$$

Since R is symmetric and has each diagonal element equal to one, its redundant and nonrandom elements can be eliminated by using $\tilde{v}(R)$. Since $\tilde{v}(R) = \tilde{L}_m \text{vec}(R)$, where \tilde{L}_m is the matrix discussed in Section 8.7, we find that the asymptotic covariance matrix of $\tilde{v}(R)$ is given by

$$\text{var}\{\tilde{v}(R)\} = \frac{2}{n-1}\tilde{L}_m N_m \Theta N_m \tilde{L}_m'.$$

Note that the Hadamard product and its associated properties can be useful in analyses involving the manipulation of Θ because

$$\begin{aligned} \Theta &= P \otimes P - (I_m \otimes P)\Lambda_m(P \otimes P) - (P \otimes P)\Lambda_m(I_m \otimes P) \\ &\quad + (I_m \otimes P)\Lambda_m(P \otimes P)\Lambda_m(I_m \otimes P), \end{aligned}$$

and the last term on the right-hand side of this equation can be expressed as

$$(I_m \otimes P)\Lambda_m(P \otimes P)\Lambda_m(I_m \otimes P) = (I_m \otimes P)\Psi_m'(P \odot P)\Psi_m(I_m \otimes P).$$

PROBLEMS

- 11.1** We saw in the proof of Theorem 11.1 that if A is an $m \times m$ idempotent matrix, then $\text{rank}(A) + \text{rank}(I_m - A) = m$. Prove the converse; that is, show that if A is an $m \times m$ matrix satisfying $\text{rank}(A) + \text{rank}(I_m - A) = m$, then A is idempotent.

- 11.2** Suppose that A is an $m \times m$ idempotent matrix. Show that each of the following matrices is also idempotent:
- (a) A' .
 - (b) BAB^{-1} , where B is any $m \times m$ nonsingular matrix.
 - (c) A^n , where n is a positive integer.
- 11.3** Show that if A is an $m \times m$ symmetric idempotent matrix having rank r , then $A = PP'$ for some $m \times r$ matrix satisfying $P'P = I_r$.
- 11.4** Let A be an $m \times n$ matrix. Show that each of the following matrices is idempotent:
- (a) AA^- .
 - (b) A^-A .
 - (c) $A(A'A)^-A'$.
- 11.5** Let A and B be $m \times m$ symmetric idempotent matrices. Show that if the column spaces of A and B are the same, then $A = B$.
- 11.6** Determine the class of $m \times 1$ vectors $\{x\}$, for which xx' is idempotent.
- 11.7** Determine the values of the scalars a , b , and c for which each of the following is an idempotent matrix.
- (a) $a\mathbf{1}_m\mathbf{1}_m'$.
 - (b) $bI_m + c\mathbf{1}_m\mathbf{1}_m'$.
- 11.8** Let A be an $m \times n$ matrix with $\text{rank}(A) = m$. Show that $A'(AA')^{-1}A$ is symmetric, idempotent, and find its rank.
- 11.9** Let A and B be $m \times m$ matrices. Show that if B is nonsingular and AB is idempotent, then BA is also idempotent.
- 11.10** Let A be an $m \times m$ symmetric idempotent matrix of rank r . Show that if B is an $m \times r$ matrix of rank r satisfying $AB = B$, then $A = B(B'B)^{-1}B'$.
- 11.11** Show that if A is an $m \times m$ matrix and $A^2 = cA$ for some scalar c , then

$$\text{tr}(A) = c \text{rank}(A).$$

- 11.12** Let A be an $m \times m$ symmetric idempotent matrix and B be an $m \times m$ nonnegative definite matrix. Show that if $I_m - A - B$ is nonnegative definite, then $AB = BA = (0)$.
- 11.13** Let A be an $m \times m$ symmetric idempotent matrix and B be an $m \times m$ matrix.
- (a) Show that if $AB = B$, then $A - BB^+$ is symmetric idempotent with rank equal to $\text{rank}(A) - \text{rank}(B)$.
 - (b) Show that if $AB = (0)$ and $\text{rank}(A) + \text{rank}(B) = m$, then $A = I_m - BB^+$.
- 11.14** Give an example of a collection of matrices A_1, \dots, A_k that satisfies conditions (a) and (d) of Corollary 11.8.1, but it does not satisfy conditions (b) and (c). Similarly, find a collection of matrices that satisfies conditions (c) and (d), but it does not satisfy conditions (a) and (b).

11.15 Prove Theorem 11.12.

11.16 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is a positive definite matrix, and let A be an $m \times m$ symmetric matrix.

(a) Show that the moment generating function of $y = \mathbf{x}'A\mathbf{x}$ can be expressed as

$$m_y(t) = |I_m - 2tA\Omega|^{-1/2} \exp \left\{ -\frac{1}{2}\boldsymbol{\mu}'[I_m - (I_m - 2tA\Omega)^{-1}]\Omega^{-1}\boldsymbol{\mu} \right\}.$$

(b) If $w \sim \chi_r^2(\frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu})$, then it can be shown that

$$m_w(t) = (1 - 2t)^{-r/2} \exp \left\{ -\frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu}[1 - (1 - 2t)^{-1}] \right\}.$$

Use this result and the moment generating function from (a) to show that the sufficient condition given in Theorem 11.12 is necessary as well. Do this by equating the two moment generating functions at $\boldsymbol{\mu} = \mathbf{0}$ and using the resulting equation to show that $A\Omega$ must be idempotent of rank r .

11.17 Let A be an $m \times m$ symmetric matrix with $r = \text{rank}(A)$, and suppose that $\mathbf{x} \sim N_m(\mathbf{0}, I_m)$. Show that the distribution of $\mathbf{x}'A\mathbf{x}$ can be expressed as a linear combination of r independent chi-squared random variables, each with one degree of freedom. What are the coefficients in this linear combination when A is idempotent?

11.18 Extend the result of Problem 11.17 to the situation in which $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, where Ω is nonnegative definite; that is, show that if A is a symmetric matrix, then $\mathbf{x}'A\mathbf{x}$ can be expressed as a linear combination of independent chi-squared random variables each having one degree of freedom. How many chi-squared random variables are in this linear combination?

11.19 Let x_1, \dots, x_n be a random sample from a normal distribution with mean μ and variance σ^2 , and let \bar{x} be the sample mean. Write

$$t = \frac{n(\bar{x} - \mu)^2}{\sigma^2}$$

as a quadratic form in the vector $(\mathbf{x} - \mu\mathbf{1}_n)$, where $\mathbf{x} = (x_1, \dots, x_n)'$. What is the distribution of t ?

11.20 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite, and suppose that A and B are $m \times m$ symmetric matrices. Show that the sufficient condition given in Theorem 11.14 is also necessary. That is, show that if $\mathbf{x}'A\mathbf{x}$ and $\mathbf{x}'B\mathbf{x}$ are independently distributed, then $A\Omega B = (0)$.

11.21 Suppose that $\mathbf{x} \sim N_n(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite. Partition \mathbf{x} , $\boldsymbol{\mu}$, and Ω as

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega'_{12} & \Omega_{22} \end{bmatrix},$$

where \mathbf{x}_1 is $r \times 1$ and \mathbf{x}_2 is $(n - r) \times 1$. Show that

- (a) $t_1 = (\mathbf{x}_1 - \boldsymbol{\mu}_1)' \Omega_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \sim \chi_r^2$,
 (b) $t_2 = (\mathbf{x} - \boldsymbol{\mu})' \Omega^{-1} (\mathbf{x} - \boldsymbol{\mu}) - (\mathbf{x}_1 - \boldsymbol{\mu}_1)' \Omega_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \sim \chi_{n-r}^2$,
 (c) t_1 and t_2 are independently distributed.
- 11.22** Suppose that $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ and the $m \times m$ matrices A_1 and A_2 are such that $t_1 = \mathbf{x}' A_1 \mathbf{x} \sim \chi_{d_1}^2$ and $t_2 = \mathbf{x}' A_2 \mathbf{x} \sim \chi_{d_2}^2$, independently. Consequently, $t = (t_1/d_1)/(t_2/d_2)$ has the F distribution with d_1 and d_2 degrees of freedom. Show that if \mathbf{y} has an elliptical distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix Ω , then $w = (w_1/d_1)/(w_2/d_2)$ has this same F distribution, where $w_1 = \mathbf{y}' A_1 \mathbf{y}$ and $w_2 = \mathbf{y}' A_2 \mathbf{y}$.
- 11.23** Prove Theorem 11.15.
- 11.24** Pearson's chi-squared statistic is given by

$$t = \sum_{i=1}^m \frac{(nx_i - n\mu_i)^2}{n\mu_i},$$

where n is a positive integer, the x_i 's are random variables, and the μ_i 's are nonnegative constants satisfying $\mu_1 + \cdots + \mu_m = 1$. Let $\mathbf{x} = (x_1, \dots, x_m)'$, $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)'$, and $\Omega = D - \boldsymbol{\mu}\boldsymbol{\mu}'$, where $D = \text{diag}(\mu_1, \dots, \mu_m)$.

- (a) Show that Ω is a singular matrix.
 (b) Show that if $\sqrt{n}(\mathbf{x} - \boldsymbol{\mu}) \sim N_m(\mathbf{0}, \Omega)$, then $t \sim \chi_{m-1}^2$.
- 11.25** Suppose that $\mathbf{x} \sim N_4(\mathbf{0}, I_4)$, and consider the three functions of the components of \mathbf{x} given by

$$\begin{aligned} t_1 &= \frac{1}{4}(x_1 + x_2 + x_3 + x_4)^2 + \frac{1}{2}(x_1 - x_2)^2, \\ t_2 &= \frac{1}{12}(x_1 + x_2 + x_3 - 3x_4)^2, \\ t_3 &= (x_1 + x_2 - 2x_3)^2 + (x_3 - x_4)^2. \end{aligned}$$

- (a) Write t_1 , t_2 , and t_3 as quadratic forms in \mathbf{x} .
 (b) Which of these statistics have chi-squared distributions?
 (c) Which of the pairs t_1 and t_2 , t_1 and t_3 , and t_2 and t_3 are independently distributed?
- 11.26** Suppose that $\mathbf{x} \sim N_4(\boldsymbol{\mu}, \Omega)$, where $\boldsymbol{\mu} = (1, -1, 1, -1)'$ and $\Omega = I_4 + \mathbf{1}_4 \mathbf{1}_4'$. Define

$$\begin{aligned} t_1 &= \frac{1}{2}(x_1 - x_2)^2 + \frac{1}{2}(x_3 - x_4)^2, \\ t_2 &= \frac{1}{2}(x_1 + x_2)^2 + \frac{1}{2}(x_3 + x_4)^2. \end{aligned}$$

(a) Does t_1 or t_2 have a chi-squared distribution? If so, identify the parameters of the distribution.

(b) Are t_1 and t_2 independently distributed?

11.27 Prove Theorem 11.16.

11.28 Prove Theorem 11.17.

11.29 The purpose of this exercise is to generalize the results of Example 11.5 to a test of the hypothesis that $H\beta = c$, where H is an $m_2 \times m$ matrix having rank m_2 and c is an $m_2 \times 1$ vector; Example 11.5 dealt with the special case in which $H = \begin{pmatrix} (0) & I_{m_2} \end{pmatrix}$ and $c = 0$. Let G be an $(m - m_2) \times m$ matrix having rank $m - m_2$ and satisfying $HG' = (0)$. Show that the reduced model may be written as

$$y_* = X_*\beta_* + \epsilon,$$

where $y_* = y - XH'(HH')^{-1}c$, $X_* = XG'(GG')^{-1}$, and $\beta_* = G\beta$. Use the sum of squared errors for this reduced model and the sum of squared errors for the complete model to construct the appropriate F statistic.

11.30 Suppose that $x \sim N_m(0, \Omega)$, where $r = \text{rank}(\Omega) < m$. If T is any $m \times r$ matrix satisfying $TT' = \Omega$, and $z \sim N_r(0, I_r)$, then x is distributed the same as Tz . Use this to show that the formulas given in Theorem 11.22 for positive definite Ω also hold when Ω is positive semidefinite.

11.31 Let $z \sim N_m(0, I_m)$. Use the fact that the first six moments of the standard normal distribution are 0, 1, 0, 3, 0, and 15 to show that

$$\begin{aligned} E(zz' \otimes zz' \otimes zz') &= I_{m^3} + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (I_m \otimes T_{ij} \otimes T_{ij} \\ &\quad + T_{ij} \otimes I_m \otimes T_{ij} + T_{ij} \otimes T_{ij} \otimes I_m) \\ &\quad + \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m (T_{ij} \otimes T_{ik} \otimes T_{jk}), \end{aligned}$$

where $T_{ij} = E_{ij} + E_{ji}$.

11.32 Suppose that $z \sim N_m(0, I_m)$.

(a) Show that

$$E(zz' \otimes zz') = N_m\{2I_{m^2} + \text{vec}(I_m)\text{vec}(I_m)'\}N_m.$$

(b) Let Δ be the matrix defined in Problem 8.61. Show that the sixth-order moment matrix given in Problem 11.31 can be more compactly expressed as

$$E(zz' \otimes zz' \otimes zz') = \Delta\{6I_{m^3} + 9I_m \otimes \text{vec}(I_m)\text{vec}(I_m)'\}\Delta.$$

Expressions for higher order moment matrices of z , such as $E(zz' \otimes zz' \otimes zz' \otimes zz')$, can be found in Schott (2003).

11.33 Suppose that \mathbf{y} has an elliptical distribution with mean vector $\mathbf{0}$, covariance matrix Ω , and finite fourth moments.

(a) Show that for some constant c ,

$$E(\mathbf{y}\mathbf{y}' \otimes \mathbf{y}\mathbf{y}') = c\{2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\text{vec}(\Omega)'\}.$$

(b) Use the expression given in (a) to show that

$$c = \frac{E(y_i^4)}{3\{E(y_i^2)\}^2},$$

regardless of the choice of i .

(c) Show that if S is the sample covariance matrix computed from a random sample of size n from this elliptical distribution, then

$$\text{var}\{\text{vec}(S)\} \approx \frac{1}{(n-1)}\{2cN_m(\Omega \otimes \Omega) + (c-1)\text{vec}(\Omega)\text{vec}(\Omega)'\},$$

for large n .

(d) If R is the sample correlation matrix, show that the first-order approximation

$$\text{var}\{\text{vec}(R)\} \approx \frac{2c}{n-1}N_m\Theta N_m$$

holds, where Θ is as defined in Example 11.11.

11.34 Suppose that \mathbf{u} is uniformly distributed on the m -dimensional unit sphere.

(a) Show that

$$E(\mathbf{u} \otimes \mathbf{u}) = m^{-1}\text{vec}(I_m).$$

(b) Show that

$$E(\mathbf{u}\mathbf{u}' \otimes \mathbf{u}\mathbf{u}') = \{m(m+2)\}^{-1}\{2N_m + \text{vec}(I_m)\text{vec}(I_m)'\}.$$

11.35 Let A , B , and C be $m \times m$ symmetric matrices, and suppose that $\mathbf{x} \sim N_m(\mathbf{0}, I_m)$.

(a) Show that

$$E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}) = \text{tr}\{(A \otimes B \otimes C)E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}')\}.$$

(b) Use part (a) and the result of Problem 11.31 to derive the formula given in Theorem 11.24 for $E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x})$.

11.36 Let $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite.

(a) Using Theorem 11.21, show that

$$\text{var}(\mathbf{x} \otimes \mathbf{x}) = 2N_m(\Omega \otimes \Omega + \Omega \otimes \boldsymbol{\mu}\boldsymbol{\mu}' + \boldsymbol{\mu}\boldsymbol{\mu}' \otimes \Omega).$$

(b) Show that the matrix $(\Omega \otimes \Omega + \Omega \otimes \boldsymbol{\mu}\boldsymbol{\mu}' + \boldsymbol{\mu}\boldsymbol{\mu}' \otimes \Omega)$ is nonsingular.

(c) Determine the eigenvalues of N_m . Use these along with part (b) to show that $\text{rank}\{\text{var}(\mathbf{x} \otimes \mathbf{x})\} = m(m+1)/2$.

11.37 Suppose that the $m \times 1$ vector \mathbf{x} and the $n \times 1$ vector \mathbf{y} are independently distributed with $E(\mathbf{x}) = \boldsymbol{\mu}_1$, $E(\mathbf{y}) = \boldsymbol{\mu}_2$, $E(\mathbf{x}\mathbf{x}') = V_1$, and $E(\mathbf{y}\mathbf{y}') = V_2$. Show that

(a) $E(\mathbf{x}\mathbf{y}' \otimes \mathbf{x}\mathbf{y}') = \text{vec}(V_1)\{\text{vec}(V_2)\}'$,

(b) $E(\mathbf{x}\mathbf{y}' \otimes \mathbf{y}\mathbf{x}') = (V_1 \otimes V_2)K_{mn} = K_{mn}(V_2 \otimes V_1)$,

(c) $E(\mathbf{x} \otimes \mathbf{x} \otimes \mathbf{y} \otimes \mathbf{y}) = \text{vec}(V_1) \otimes \text{vec}(V_2)$,

(d) $E(\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{x} \otimes \mathbf{y}) = (I_m \otimes K_{nm} \otimes I_n)\{\text{vec}(V_1) \otimes \text{vec}(V_2)\}$,

(e) $\text{var}(\mathbf{x} \otimes \mathbf{y}) = V_1 \otimes V_2 - \boldsymbol{\mu}_1\boldsymbol{\mu}_1' \otimes \boldsymbol{\mu}_2\boldsymbol{\mu}_2'$.

11.38 Let A , B , and C be $m \times m$ symmetric matrices, and let \mathbf{a} and \mathbf{b} be $m \times 1$ vectors of constants. If $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$, show that

(a) $E(\mathbf{x}'A\mathbf{a}\mathbf{x}'B\mathbf{b}) = \mathbf{a}'A\Omega B\mathbf{b}$,

(b) $E(\mathbf{x}'A\mathbf{a}\mathbf{x}'B\mathbf{b}\mathbf{x}'C\mathbf{x}) = \mathbf{a}'A\Omega B\mathbf{b}\text{tr}(\Omega C) + 2\mathbf{a}'A\Omega C\Omega B\mathbf{b}$.

11.39 Suppose that $\mathbf{x} \sim N_4(\boldsymbol{\mu}, \Omega)$, where $\boldsymbol{\mu} = \mathbf{1}_4$ and $\Omega = 4I_4 + \mathbf{1}_4\mathbf{1}_4'$. Let the random variables t_1 and t_2 be defined by

$$t_1 = (x_1 + x_2 - 2x_3)^2 + (x_3 - x_4)^2,$$

$$t_2 = (x_1 - x_2 - x_3)^2 + (x_1 + x_2 - x_4)^2.$$

Use Theorem 11.23 to find

(a) $\text{var}(t_1)$,

(b) $\text{var}(t_2)$,

(c) $\text{cov}(t_1, t_2)$.

11.40 Verify the formulas given at the end of Example 11.7 for $E(t_1)$ and $\text{var}(t_1)$.

11.41 Suppose that $V_1 \sim W_m(\Omega, n_1)$ and $V_2 \sim W_m(\Omega, n_2)$ are independently distributed. Show that $V_1 + V_2 \sim W_m(\Omega, n_1 + n_2)$.

11.42 Suppose that $V \sim W_m(\Omega, n, \Phi)$ and \mathbf{a} is a nonnull $m \times 1$ vector of constants. Show that $\mathbf{a}'V\mathbf{a}/\mathbf{a}'\Omega\mathbf{a} \sim \chi_n^2(\lambda)$, where $\lambda = \mathbf{a}'\Phi\mathbf{a}/\mathbf{a}'\Omega\mathbf{a}$.

11.43 Consider the Wishart matrix V given in Theorem 11.27.

(a) Show that $V_{11} - V_{12}V_{22}^{-1}V_{12}'$ is independently distributed of V_{12} and V_{22} .

(b) Show that the conditional distribution of V_{12} given V_{22} is multivariate normal; in particular, show that given V_{22} ,

$$\text{vec}(V_{12}) \sim N_{m_1m_2}(\text{vec}(\Omega_{12}\Omega_{22}^{-1}V_{22}), V_{22} \otimes (\Omega_{11} - \Omega_{12}\Omega_{22}^{-1}\Omega_{12}')).$$

- 11.44** Let $V \sim W_m(\Omega, n)$, where Ω is a positive definite matrix, and let V_k and Ω_k be the leading $k \times k$ principal submatrices of V and Ω ; that is, V_k is the matrix obtained by deleting the last $m - k$ rows and columns of V , and similarly for Ω_k . Show that if we define $|V_0| = 1$ and $|\Omega_0| = 1$, then

$$\frac{|V_k|}{|V_{k-1}|} \frac{|\Omega_{k-1}|}{|\Omega_k|} \sim \chi_{n-k+1}^2$$

for $k = 1, \dots, m$.

- 11.45** Suppose that $V \sim W_m(\Omega, n)$, where Ω is positive definite. Use Theorem 11.27 to show that if A is a $p \times m$ matrix of rank p , then $(AV^{-1}A')^{-1} \sim W_p((A\Omega^{-1}A')^{-1}, n - m + p)$.
- 11.46** Suppose that $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$, where Ω is positive definite. Partition \mathbf{x} as $\mathbf{x} = (\mathbf{x}'_1, \mathbf{x}'_2)'$, where \mathbf{x}_1 is $m_1 \times 1$ and \mathbf{x}_2 is $m_2 \times 1$. Similarly, $\boldsymbol{\mu}$ and Ω are partitioned as

$$\boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega'_{12} & \Omega_{22} \end{bmatrix}.$$

- (a) Show that $E(\mathbf{x}_1 \otimes \mathbf{x}_2) = \text{vec}(\Omega'_{12}) + \boldsymbol{\mu}_1 \otimes \boldsymbol{\mu}_2$.
- (b) Show that

$$\begin{aligned} \text{var}(\mathbf{x}_1 \otimes \mathbf{x}_2) &= \Omega_{11} \otimes \Omega_{22} + \Omega_{11} \otimes \boldsymbol{\mu}_2 \boldsymbol{\mu}'_2 + \boldsymbol{\mu}_1 \boldsymbol{\mu}'_1 \otimes \Omega_{22} \\ &\quad + K_{m_1 m_2}(\Omega'_{12} \otimes \Omega_{12} + \Omega'_{12} \otimes \boldsymbol{\mu}_1 \boldsymbol{\mu}'_2 \\ &\quad + \boldsymbol{\mu}_2 \boldsymbol{\mu}'_1 \otimes \Omega_{12}). \end{aligned}$$

- 11.47** Suppose that the columns of $X' = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ are independently distributed with $\mathbf{x}_i \sim N_m(\boldsymbol{\mu}_i, \Omega)$. Let A be an $n \times n$ symmetric matrix, and let $M' = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_n)$. Use the spectral decomposition of A to show that

- (a) $E(X'AX) = \text{tr}(A)\Omega + M'AM$,
- (b) $\text{var}\{\text{vec}(X'AX)\} = 2N_m\{\text{tr}(A^2)(\Omega \otimes \Omega) + \Omega \otimes M'A^2M + M'A^2M \otimes \Omega\}$.

- 11.48** Let A and B be $m \times n$ matrices of constants, whereas $\mathbf{x} \sim N_n(\boldsymbol{\mu}, \Omega)$. Note that $(A\mathbf{x} \odot B\mathbf{x}) = \Psi_m(A\mathbf{x} \otimes B\mathbf{x})$, where Ψ_m is the matrix defined in Section 8.5.

- (a) Show that

$$E(A\mathbf{x} \odot B\mathbf{x}) = D_{B\Omega A'}\mathbf{1}_m + A\boldsymbol{\mu} \odot B\boldsymbol{\mu},$$

where $D_{B\Omega A'}$ is the diagonal matrix with diagonal elements equal to those of $B\Omega A'$.

(b) Show that

$$\begin{aligned} \text{var}(A\mathbf{x} \odot B\mathbf{x}) &= A(\Omega + \boldsymbol{\mu}\boldsymbol{\mu}')A' \odot B(\Omega + \boldsymbol{\mu}\boldsymbol{\mu}')B' \\ &\quad + B(\Omega + \boldsymbol{\mu}\boldsymbol{\mu}')A' \odot A(\Omega + \boldsymbol{\mu}\boldsymbol{\mu}')B' \\ &\quad - A\boldsymbol{\mu}\boldsymbol{\mu}'A' \odot B\boldsymbol{\mu}\boldsymbol{\mu}'B' - B\boldsymbol{\mu}\boldsymbol{\mu}'A' \odot A\boldsymbol{\mu}\boldsymbol{\mu}'B'. \end{aligned}$$

For some applications of these results as well as generalizations, see Hyn-dman and Wand (1997), Neudecker and Liu (2001), Neudecker, et al. (1995a), and Neudecker, et al. (1995b).

11.49 Suppose that the $m \times 1$ vectors $\{\mathbf{y}_{ij}, 1 \leq i \leq k, 1 \leq j \leq n_i\}$ are independently distributed with $\mathbf{y}_{ij} \sim N_m(\boldsymbol{\mu}_i, \Omega)$. A multivariate analysis of variance uses the matrices (Example 3.16)

$$B = \sum_{i=1}^k n_i(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})(\bar{\mathbf{y}}_i - \bar{\mathbf{y}})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)(\mathbf{y}_{ij} - \bar{\mathbf{y}}_i)',$$

where

$$\bar{\mathbf{y}}_i = \sum_{j=1}^{n_i} \frac{\mathbf{y}_{ij}}{n_i}, \quad \bar{\mathbf{y}} = \sum_{i=1}^k \frac{n_i \bar{\mathbf{y}}_i}{n}, \quad n = \sum_{i=1}^k n_i.$$

Use Theorem 11.25 to show that W and B are independently distributed, $W \sim W_m(\Omega, w)$, and $B \sim W_m(\Omega, b, \Phi)$, where $w = n - k$, $b = k - 1$, and

$$\Phi = \frac{1}{2} \sum_{i=1}^k n_i(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})', \quad \bar{\boldsymbol{\mu}} = \sum_{i=1}^k \frac{n_i \boldsymbol{\mu}_i}{n}.$$

11.50 Let $X' = (\mathbf{x}_1, \dots, \mathbf{x}_n)$ be an $m \times n$ matrix, where $\mathbf{x}_1, \dots, \mathbf{x}_n$ are independent and $\mathbf{x}_i \sim N_m(\mathbf{0}, \Omega)$ for each i . Show that

$$\begin{aligned} E(X \otimes X \otimes X \otimes X) &= \{\text{vec}(I_n) \otimes \text{vec}(I_n)\} \{\text{vec}(\Omega) \otimes \text{vec}(\Omega)\}' \\ &\quad + \text{vec}(I_n \otimes I_n) \{\text{vec}(\Omega \otimes \Omega)\}' + \text{vec}(K_{nn}) \\ &\quad \times [\text{vec}\{K_{mm}(\Omega \otimes \Omega)\}]'. \end{aligned}$$

11.51 Let the columns of the $m \times n$ matrix X' be independently and identically distributed as $N_m(\mathbf{0}, \Omega)$. Suppose the $n \times m$ matrix M and the $n \times n$ matrix A contain constants, and define $V = (X + M)'A(X + M)$. Show that

$$\begin{aligned} \text{var}\{\text{vec}(V)\} &= \{\text{tr}(A'A)\}(\Omega \otimes \Omega) + \{\text{tr}(A^2)\}K_{mm}(\Omega \otimes \Omega) \\ &\quad + M'A'AM \otimes \Omega + \Omega \otimes M'AA'M \\ &\quad + K_{mm}(M'A^2M \otimes \Omega) + K_{mm}(\Omega \otimes M'A^2M)'. \end{aligned}$$

- 11.52** Suppose that the smallest eigenvalue of the $m \times m$ covariance matrix Ω has multiplicity r , and let P denote the eigenprojection of Ω corresponding to this smallest eigenvalue. Let S be the sample covariance matrix computed from a random sample of size n , and define $A = S - \Omega$. Then

$$U = \frac{r^{-1} \sum_{i=m-r+1}^m \lambda_i^2(S)}{\{r^{-1} \sum_{i=m-r+1}^m \lambda_i(S)\}^2} - 1,$$

where $\lambda_1(S) \geq \cdots \geq \lambda_m(S)$ are the eigenvalues of S , has the second-order approximation formula in A (see Problem 9.26) given by

$$U \approx r^{-1}(\text{tr}(APAP) - r^{-1}\{\text{tr}(AP)\}^2).$$

Use this approximation to show that, when sampling from a normal population, $nrU/2$ can be approximated by the chi-squared distribution with $r(r+1)/2 - 1$ degrees of freedom.

- 11.53** Use the results of Problem 8.48(e) and Problem 8.60 to show that

$$\left(I_{m^2} - \frac{1}{2} \{ (I_m \otimes P) + (P \otimes I_m) \} \Lambda_m \right) N_m = N_m \{ I_{m^2} - (I_m \otimes P) \Lambda_m \}$$

thereby verifying the simplified formula for $\text{var}\{\text{vec}(R)\}$ given in (11.18).

- 11.54** Let S be the $m \times m$ sample covariance matrix computed from a sample of size n from a normal population with covariance matrix Ω . Denote the eigenvalues and normalized eigenvectors of Ω by $x_1 \geq \cdots \geq x_m$ and $\mathbf{q}_1, \dots, \mathbf{q}_m$, and those of S by $\lambda_1 \geq \cdots \geq \lambda_m$ and $\gamma_1, \dots, \gamma_m$. Suppose that $x_k > x_{k+1} = \cdots = x_m$, and consider the eigenprojection, $P = \sum_{i=k+1}^m \mathbf{q}_i \mathbf{q}_i'$, associated with the eigenvalues x_{k+1}, \dots, x_m . An estimate of this eigenprojection is given by $\hat{P} = \sum_{i=k+1}^m \gamma_i \gamma_i'$. Use Theorem 9.7 and the large sample distribution of S discussed in Example 11.9 to show that for large n , $\text{vec}(\hat{P}) \sim N_{m^2}(\text{vec}(P), 2N_m \Psi / (n-1))$, approximately, where the matrix Ψ is given by

$$\Psi = \sum_{i=1}^k \sum_{j=k+1}^m \frac{x_i x_j}{(x_i - x_j)^2} (\mathbf{q}_i \mathbf{q}_i' \otimes \mathbf{q}_j \mathbf{q}_j' + \mathbf{q}_j \mathbf{q}_j' \otimes \mathbf{q}_i \mathbf{q}_i').$$

REFERENCES

- Agaian, S. S. (1985). *Hadamard Matrices and Their Applications*. Springer-Verlag, Berlin.
- Anderson, T. W. (1955). The integral of a symmetric unimodal function over a symmetric convex set and some probability inequalities. *Proceedings of the American Mathematical Society*, **6**, 170–176.
- Anderson, T. W. (1996). Some inequalities for symmetric convex sets with applications. *Annals of Statistics*, **24**, 753–762.
- Anderson, T. W. and Das Gupta, S. (1963). Some inequalities on characteristic roots of matrices. *Biometrika*, **50**, 522–524.
- Ando, T. (1989). Majorization, doubly stochastic matrices, and comparison of eigenvalues, *Linear Algebra and Its Applications*, **118**, 163–248.
- Ando, T. (1994). Majorizations and inequalities in matrix theory. *Linear Algebra and Its Applications*, **199**, 17–67.
- Andrilli, S. and Hecker, D. (2010). *Elementary Linear Algebra*, 4th ed., Academic Press, New York.
- Barndorff-Nielsen, O. E. and Cox, D. R. (1994). *Inferences and Asymptotics*. Chapman and Hall, London.
- Bartlett, M. S. (1937). Properties of sufficiency and statistical tests. *Proceedings of the Royal Society of London, Ser. A*, **160**, 268–282.
- Bartlett, M. S. (1947). Multivariate analysis. *Journal of the Royal Statistical Society Supplement, Ser. B*, **9**, 176–197.
- Basilevsky, A. (1983). *Applied Matrix Algebra in the Statistical Sciences*. North-Holland, New York.

- Bellman, R. (1970). *Introduction to Matrix Analysis*. McGraw-Hill, New York.
- Ben-Israel, A. (1966). A note on an iterative method for generalized inversion of matrices. *Mathematics of Computation*, **20**, 439–440.
- Ben-Israel, A. and Greville, T. N. E. (2003). *Generalized Inverses: Theory and Applications*, 2nd ed. Springer-Verlag, New York.
- Berkovitz, L. D. (2002). *Convexity and Optimization in R^n* . John Wiley, New York.
- Berman, A. and Plemmons, R. J. (1994). *Nonnegative Matrices in the Mathematical Sciences*. Society for Industrial and Applied Mathematics, Singapore.
- Berman, A. and Shaked-Monderer, N. (2003). *Completely Positive Matrices*. World Scientific, Singapore.
- Bhatia, R. (1997). *Matrix Analysis*. Springer-Verlag, New York.
- Bhattacharya, R. N. and Waymire, E. C. (2009). *Stochastic Processes with Applications*. Society for Industrial and Applied Mathematics, Philadelphia.
- Boullion, T. L. and Odell, P. L. (1971). *Generalized Inverse Matrices*. John Wiley, New York.
- Campbell, S. L. and Meyer, C. D. (1979). *Generalized Inverses of Linear Transformations*. Pitman, London.
- Casella, G. and Berger, R. L. (2002). *Statistical Inference*. Duxbury, Pacific Grove, CA.
- Cline, R. E. (1964a). Note on the generalized inverse of the product of matrices. *SIAM Review*, **6**, 57–58.
- Cline, R. E. (1964b). Representations for the generalized inverse of a partitioned matrix. *SIAM Journal of Applied Mathematics*, **12**, 588–600.
- Cline, R. E. (1965). Representations for the generalized inverse of sums of matrices. *SIAM Journal of Numerical Analysis*, **2**, 99–114.
- 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.
- Davis, P. J. (1994). *Circulant Matrices*, 2nd ed. AMS Chelsea Publishing, Providence, RI.
- Duff, I. S., Erisman, A. M., and Reid, J. K. (1986). *Direct Methods for Sparse Matrices*, Oxford University Press, Oxford.
- Dümbgen, L. (1995). A simple proof and refinement of Wielandt's eigenvalue inequality. *Statistics & Probability Letters*, **25**, 113–115.
- Eaton, M. L. and Tyler, D. E. (1991). On Wielandt's inequality and its application to the asymptotic distribution of the eigenvalues of a random symmetric matrix. *Annals of Statistics*, **19**, 260–271.
- Elsner, L. (1982). On the variation of the spectra of matrices. *Linear Algebra and Its Applications*, **47**, 127–138.
- Eubank, R. L. and Webster, J. T. (1985). The singular-value decomposition as a tool for solving estimability problems. *American Statistician*, **39**, 64–66.
- Fan, K. (1949). On a theorem of Weyl concerning eigenvalues of linear transformations, I. *Proceedings of the National Academy of Sciences of the USA*, **35**, 652–655.
- Fan, K. (1950). On a theorem of Weyl concerning eigenvalues of linear transformations, II. *Proceedings of the National Academy of Sciences of the USA*, **36**, 31–35.
- Fang, K. T., Kotz, S., and Ng, K. W. (1990). *Symmetric Multivariate and Related Distributions*. Chapman and Hall, London.

- Ferguson, T. S. (1967). *Mathematical Statistics: A Decision Theoretic Approach*. Academic Press, New York.
- Gantmacher, F. R. (1959). *The Theory of Matrices*, Volumes I and II. Chelsea, New York.
- Golub, G. H. and Van Loan, C. F. (2013). *Matrix Computations*, 4th ed. Johns Hopkins University Press, Baltimore.
- Graybill, F. A. (1983). *Matrices With Applications in Statistics*, 2nd ed. Wadsworth, Belmont, CA.
- Grenander, U. and Szego, G. (1984). *Toeplitz Forms and Their Applications*. Chelsea, New York.
- Greville, T. N. E. (1960). Some applications of the pseudoinverse of a matrix. *SIAM Review*, **2**, 15–22.
- Greville, T. N. E. (1966). Note on the generalized inverse of a matrix product. *SIAM Review*, **8**, 518–521.
- Gross, J. (2000). The Moore–Penrose inverse of a partitioned nonnegative definite matrix. *Linear Algebra and its Applications*, **321**, 113–121.
- Gustafson, K. (1972). Antieigenvalue inequalities in operator theory. In *Inequalities III* (O. Shisha, ed.), 115–119. Academic Press, New York.
- Gustafson, K. (2006). The trigonometry of matrix statistics. *International Statistical Review*, **74**, 187–202.
- Hageman, L. A. and Young, D. M. (1981). *Applied Iterative Methods*. Academic Press, New York.
- Hardy, G. H., Littlewood, J. E., and Pólya, G. (1952). *Inequalities*, 2nd ed. Cambridge University Press, Cambridge.
- Harville, D. A. (1997). *Matrix Algebra From a Statistician's Perspective*. Springer, New York.
- Healy, M. J. R. (1986). *Matrices for Statistics*. Clarendon Press, Oxford.
- Hedayat, A. and Wallis, W. D. (1978). Hadamard matrices and their applications. *Annals of Statistics*, **6**, 1184–1238.
- Heinig, G. and Rost, K. (1984). *Algebraic Methods for Toeplitz-like Matrices and Operators*. Birkhäuser, Basel.
- Henderson, H. V. and Searle, S. R. (1979). Vec and vech operators for matrices with some uses in Jacobians and multivariate statistics. *Canadian Journal of Statistics*, **7**, 65–81.
- Hinch, E. J. (1991). *Perturbation Methods*. Cambridge University Press, Cambridge.
- Horn, R. A. and Johnson, C. R. (1985). *Matrix Analysis*. Cambridge University Press, Cambridge.
- Horn, R. A. and Johnson, C. R. (1991). *Topics in Matrix Analysis*. Cambridge University Press, Cambridge.
- Horn, R. A. and Johnson, C. R. (2013). *Matrix Analysis*, 2nd ed. Cambridge University Press, Cambridge.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, **24**, 417–441, 498–520.
- Hu, X. (2008). A three-condition characterization of the Moore–Penrose generalized inverse. *The American Statistician*, **62**, 216–218.
- Huberty, C. J and Olejnik, S. (2006). *Applied MANOVA and Discriminant Analysis*, 2nd ed. John Wiley, New York.

- Hyndman, R. J. and Wand, M. P. (1997). Nonparametric autocovariance function estimation. *Australian Journal of Statistics*, **39**, 313–324.
- Im, E. I. (1997). Narrower eigenbounds for Hadamard products. *Linear Algebra and Its Applications*, **254**, 141–144.
- Jackson, J. E. (1991). *A User's Guide to Principal Components*. John Wiley, New York.
- Johnson, N. L., Kotz, S., and Balakrishnan, N. (1997). *Discrete Multivariate Distributions*. John Wiley, New York.
- Jolliffe, I. T. (2002). *Principal Component Analysis*, 2nd ed. Springer-Verlag, New York.
- Kato, T. (1982). *A Short Introduction to Perturbation Theory for Linear Operators*. Springer-Verlag, New York.
- Kelly, P. J. and Weiss, M. L. (1979). *Geometry and Convexity*. John Wiley, New York.
- Khattree, R. (2003). Antieigenvalues and antieigenvectors in statistics. *Journal of Statistical Planning and Inference*, **114**, 131–144.
- Khuri, A. (2003). *Advanced Calculus with Applications in Statistics*, 2nd ed. John Wiley, New York.
- Krzanowski, W. J. (2000). *Principles of Multivariate Analysis: A User's Perspective*, Revised ed. Clarendon Press, Oxford.
- Kutner, M., Nachtsheim, C., Neter, J. and Li, W. (2005). *Applied Linear Statistical Models*, 5th ed. McGraw-Hill, New York.
- Lanczos, C. (1950). An iteration method for the solution of the eigenvalue problem of linear differential and integral operators. *Journal of Research of the National Bureau of Standards*, **45**, 255–282.
- Lay, S. R. (1982). *Convex Sets and Their Applications*. John Wiley, New York.
- Lidskii, V. (1950). The proper values of the sum and product of symmetric matrices. *Dokl. Akad. Nauk. SSSR*, **75**, 769–772 (in Russian). (Translated by C. D. Benster, U. S. Department of Commerce, National Bureau of Standards, Washington, D.C., N.B.S. Rep. 2248, 1953).
- Lindgren, B. W. (1993). *Statistical Theory*, 4th ed. Chapman and Hall, New York.
- Magnus, J. R. (1978). The moments of products of quadratic forms in normal variables. *Statistica Neerlandica*, **32**, 201–210.
- Magnus, J. R. (1987). A representation theorem for $(\text{tr} A^p)^{1/p}$. *Linear Algebra and Its Applications*, **95**, 127–134.
- Magnus, J. R. (1988). *Linear Structures*. Charles Griffin, London.
- Magnus, J. R. and Neudecker, H. (1979). The commutation matrix: Some properties and applications. *Annals of Statistics*, **7**, 381–394.
- Magnus, J. R. and Neudecker, H. (1988). *Matrix Differential Calculus with Applications in Statistics and Econometrics*. John Wiley, New York.
- Magnus, J. R. and Neudecker, H. (1999). *Matrix Differential Calculus with Applications in Statistics and Econometrics*, Revised ed. John Wiley, New York.
- Mandel, J. (1982). Use of the singular value decomposition in regression analysis. *American Statistician*, **36**, 15–24.
- Mardia, K. V., Kent, J. T., and Bibby, J. M. (1979). *Multivariate Analysis*. Academic Press, New York.
- Marshall, A. W., Olkin, I., and Arnold, B. C. (2011). *Inequalities: Theory of Majorization and Its Applications*, 2nd ed. Springer, New York.

- Mathai, A. M. and Provost, S. B. (1992). *Quadratic Forms in Random Variables*. Marcel Dekker, New York.
- McCullagh, P. (1987). *Tensor Methods in Statistics*. Chapman and Hall, London.
- McLachlan, G. J. (2005). *Discriminant Analysis and Statistical Pattern Recognition*. John Wiley, New York.
- Medhi, J. (2009). *Stochastic Processes*, 3rd ed. New Age Science, New Delhi.
- Miller, R. G., Jr. (1981). *Simultaneous Statistical Inference*, 2nd ed. Springer-Verlag, New York.
- Minc, H. (1988). *Nonnegative Matrices*. John Wiley, New York.
- Moore, E. H. (1920). On the reciprocal of the general algebraic matrix (Abstract). *Bulletin of the American Mathematical Society*, **26**, 394–395.
- Moore, E. H. (1935). General analysis. *Memoirs of the American Philosophical Society*, **1**, 147–209.
- Morrison, D. F. (2005). *Multivariate Statistical Methods*, 4th ed. McGraw-Hill, New York.
- Muirhead, R. J. (1982). *Aspects of Multivariate Statistical Theory*. John Wiley, New York.
- Nayfeh, A. H. (1981). *Introduction to Perturbation Techniques*. John Wiley, New York.
- Nel, D. G. (1980). On matrix differentiation in statistics. *South African Statistical Journal*, **14**, 137–193.
- Nelder, J. A. (1985). An alternative interpretation of the singular-value decomposition in regression. *American Statistician*, **39**, 63–64.
- Neudecker, H. and Liu, S. (2001). Statistical properties of the Hadamard product of random vectors. *Statistical Papers*, **42**, 529–533.
- Neudecker, H., Liu, S., and Polasek, W. (1995a). The Hadamard product and some of its applications in statistics. *Statistics*, **26**, 365–373.
- Neudecker, H., Polasek, W., and Liu, S. (1995b). The heteroskedastic linear regression model and the Hadamard product: A note. *Journal of Econometrics*, **68**, 361–366.
- Olkin, I. and Tomsy, J. L. (1981). A new class of multivariate tests based on the union-intersection principle. *Annals of Statistics*, **9**, 792–802.
- Ostrowski, A. M. (1973). *Solutions of Equations in Euclidean and Banach Spaces*. Academic Press, New York.
- Ouellette, D. V. (1981). Schur complements and statistics. *Linear Algebra and its Applications*, **36**, 187–295.
- Penrose, R. (1955). A generalized inverse for matrices. *Proceedings of the Cambridge Philosophical Society*, **51**, 406–413.
- Penrose, R. (1956). On best approximate solutions of linear matrix equations. *Proceedings of the Cambridge Philosophical Society*, **52**, 17–19.
- Perlman, M. D. (1990). T. W. Anderson's Theorem on the integral of a symmetric unimodal function over a symmetric convex set and its applications in probability and statistics. In *The Collected Papers of T. W. Anderson, 1943–1985* (George P. H. Styan, ed.), **2**, 1627–1641. John Wiley, New York.
- Pinsky, M. A. and Karlin, S. (2011). *An Introduction to Stochastic Modeling*, 4th ed. Academic Press, Burlington, MA.
- Poincaré, H. (1890). Sur les équations aux dérivées partielles de la physique mathématique. *American Journal of Mathematics*, **12**, 211–294.

- Poole, D. (2015). *Linear Algebra: A Modern Introduction*, 4th ed., Cengage Learning, Stamford, CT.
- Press, W. H., Teukolsky, S. A., Vetterline, W. T., and Flannery, B. P. (2007). *Numerical Recipes: The Art of Scientific Computing*, 3rd ed. Cambridge University Press, Cambridge.
- Pringle, R. M. and Rayner, A. A. (1971). *Generalized Inverse Matrices with Applications to Statistics*. Charles Griffin, London.
- Rao, C. R. (1973). *Linear Statistical Inference and Its Applications*, 2nd ed. John Wiley, New York.
- Rao, C. R. (2005). Antieigenvalues and antisingularvalues of a matrix and applications to problems in statistics. *Research Letters in the Information and Mathematical Sciences*, **8**, 53–76.
- Rao, C. R. and Mitra, S. K. (1971). *Generalized Inverse of Matrices and Its Applications*. John Wiley, New York.
- Rencher, A. C. and Schaafje, G. B. (2008). *Linear Models in Statistics*, 2nd ed. John Wiley, New York.
- Rockafellar, R. T. (1970). *Convex Analysis*. Princeton University Press, Princeton.
- Scheffé, H. (1953). A method for judging all contrasts in the analysis of variance. *Biometrika*, **40**, 87–104.
- Schott, J. R. (1984). Optimal bounds for the distribution of some test criteria for tests of dimensionality. *Biometrika*, **71**, 561–567.
- Schott, J. R. (2003). Kronecker product permutation matrices and their application to moment matrices of the normal distribution. *Journal of Multivariate Analysis*, **87**, 177–190.
- Searle, S. R. (1971). *Linear Models*. John Wiley, New York.
- Searle, S. R. (1982). *Matrix Algebra Useful for Statistics*. John Wiley, New York.
- Sen, A. K. and Srivastava, M. S. (1990). *Regression Analysis: Theory, Methods, and Applications*, Springer-Verlag, New York.
- Seneta, E. (2006). *Non-negative Matrices and Markov Chains*, 2nd ed. Springer, New York.
- Srivastava, M. S. and Khatri, C. G. (1979). *An Introduction to Multivariate Statistics*. North-Holland, New York.
- Stewart, G. W. (1998). *Matrix Algorithms I: Basic Decompositions*. SIAM, Philadelphia.
- Stewart, G. W. (2001). *Matrix Algorithms II: Eigensystems*. SIAM, Philadelphia.
- Styan, G. P. H. (1973). Hadamard products and multivariate statistical analysis. *Linear Algebra and Its Applications*, **6**, 217–240.
- Sugiura, N. (1976). Asymptotic expansions of the distributions of the latent roots and latent vector of the Wishart and multivariate F matrices. *Journal of Multivariate Analysis*, **6**, 500–525.
- Trenkler, G. (2000). On a generalization of the covariance matrix of the multinomial distribution. In *Innovations in Multivariate Statistical Analysis*, R.D.H. Heijmans, D.S.G. Pollock, and A. Santorra, Eds., pp. 67–73. Kluwer, Boston.
- Wielandt, H. (1955). An extremum property of sums of eigenvalues. *Proceedings of the American Mathematical Society*, **6**, 106–110.
- Xian, Y. Y. (2001). *Theory and Applications of Higher-Dimensional Hadamard Matrices*. Kluwer, Boston.
- Young, D. M. (1971). *Iterative Solution of Large Linear Systems*. Academic Press, New York.
- Zhang, F. (2005). *The Schur Complement and Its Applications*. Springer, New York.
- Zhang, F. (2011). *Matrix Theory: Basic Results and Techniques*, 2nd ed. Springer, New York.

INDEX

- accumulation point, 81
- adjoint, 9
- analysis of variance, 138, 320
 - multivariate, 493
- angle between vectors, 41, 142
- antieigenvalue, 141–144
- antieigenvector, 141–144
- arithmetic-geometric mean inequality, 452–453

- backward shift, 367
- Bartlett adjustment, 493
- basis, 49–53
 - orthonormal, 53–58
- bilinear form, 16
- block diagonal matrix, 13
- boundary point, 82

- canonical variate analysis, 119, 178–179, 493
- Cauchy–Schwarz inequality, 38, 444–446
- Cayley–Hamilton theorem, 105–106, 196
- chain rule, 388, 391
- characteristic equation, 96
- characteristic root, 96
- characteristic vector, 96
- chi-squared distribution, 21–22
 - and Moore–Penrose inverse, 211
 - and quadratic forms, 465–471
- central, 21–22
- noncentral, 22
- Cholesky decomposition, 164
- circulant matrix, 363–367
- closed set, 81
- closure, 81
- Cochran’s theorem, 462–465
- cofactor, 6
 - and determinant, 6
 - and inverse, 9
- column space, 45
- commutation matrix, 339–346
 - and Kronecker product, 342
 - and vec of a Kronecker product, 343
 - and vec operator, 341
- eigenvalues, 344
- eigenvectors, 379
- completely positive matrix, 383
- complex matrix, 18–19
- complex number, 18
 - conjugate, 18
 - Euler’s formula, 18
 - modulus, 18
 - polar coordinates, 18
 - triangle inequality, 19

- condition number, 191
- conjugate transpose, 19
- consistent equations, 247–251
- consistent estimator, 223–224
- contaminated normal distribution, 27, 34
- continuity
 - of a determinant, 222
 - of a Moore–Penrose inverse, 222–224
 - of an eigenprojection, 115
 - of an eigenvalue, 115
 - of an inverse matrix, 222
- convergence with respect to a norm, 188
- convex combination, 81
- convex function, 413–416
 - absolute minimum of, 415
 - strictly, 413
- convex hull, 81
- convex set, 80–85
- correlation, 24
 - maximum squared, 432
- correlation matrix, 24
 - eigenvalues and eigenvectors, 113–114, 331
 - nonnegative definite, 25
 - sample, 25
- Courant–Fischer min–max theorem, 120
- covariance, 23
 - of quadratic forms, 477, 480–481
- covariance matrix, 23
 - equal variances and equal correlations, 113
 - nonnegative definite, 24
 - sample, 25
- decomposition
 - Cholesky, 164
 - Jordan, 173–174
 - LU , 198
 - polar, 192
 - QR , 60, 165
 - Schur, 175–178
 - singular value, 155–162
 - spectral, 108, 111, 162–168
- density function, 20
 - multivariate, 22
- derivative, 387–389
 - of a determinant, 396, 401
 - of a matrix function, 391
 - of a Moore–Penrose inverse, 398–400
 - of a patterned matrix, 400–401
 - of a trace, 396
 - of a vector function, 390
 - of an eigenvalue, 407
 - of an eigenvector, 407
 - of an inverse, 398, 401
 - partial, 389
 - second-order partial, 390
- determinant, 5–9
 - and eigenvalues, 101
 - and trace, 452
 - continuity of, 222
 - derivative of, 396, 401
 - expansion by cofactors, 6–8
 - of a partitioned matrix, 288–296
 - of a product, 9, 445–446
 - of a sum, 102, 447, 450
- determinantal inequality, 181, 291, 336, 445–447, 450–452, 455–456
- diagonal matrix, 2
- diagonalization, 103, 169, 171–173
 - simultaneous, 136, 178–184
- differential, 388
 - of a determinant, 396
 - of a matrix function, 391
 - of a Moore–Penrose inverse, 398
 - of a trace, 396
 - of a vector function, 390
 - of an eigenvalue, 407
 - of an eigenvector, 407
 - of an inverse, 398
 - second, 389
- dimension of a vector space, 49–53
- direct sum of matrices, 323
- discriminant analysis, 40
- distance function, 39
 - Euclidean, 40, 55, 68, 166
 - Mahalanobis, 40, 68, 166
- distance in the metric of, 40
- doubly stochastic matrix, 434
- duplication matrix, 346–349
- eigenprojection, 110–112
 - continuity of, 115
- eigenspace, 98, 172
 - dimension of, 100
- eigenvalue, 95
 - and antieigenvalue, 142
 - and determinant, 101
 - and leading principal submatrix, 124
 - and majorization, 437–442
 - and rank, 104, 112, 171–173, 178
 - and trace, 101
 - asymptotic distribution of, 491
 - continuity of, 115
 - derivative of, 407
 - distinct, 98
 - extremal properties, 116–123
 - in the metric of, 136
 - monotonicity, 133
 - multiple, 98

- of a partitioned matrix, 302–307
 - of a positive definite matrix, 129
 - of a positive semidefinite matrix, 129
 - of a power, 100
 - of a product, 140–141
 - of a sum, 124–129, 133, 439–440
 - of a symmetric matrix, 106–114
 - of a transpose product, 131
 - of a triangular matrix, 99
 - of an idempotent matrix, 457
 - of an inverse matrix, 100
 - of an orthogonal matrix, 99
 - of the Schur complement, 306
 - perturbation of, 403–406, 427
 - relation to diagonal elements, 437
 - simple, 98
- eigenvector, 95
 - and antieigenvector, 142
 - asymptotic distribution of, 491
 - common, 151, 184
 - derivative of, 407
 - left, 148
 - linear independence of, 103
 - of a symmetric matrix, 107–108
 - right, 148
- elementary transformations, 14
- elimination matrix, 349–351
- elliptical distribution, 27
- estimable function, 268
- Euclidean distance function, 40, 55, 68, 166
- Euclidean inner product, 39–40
- Euclidean norm, 40, 42, 186
- Euclidean space, 40
- expected value, 20
 - of a quadratic form, 477–485
- F distribution, 22
- forward shift, 367
- Fourier matrix, 366
- Gauss–Seidel method, 274
- generalized inverse, 225–230
 - and projection matrices, 231
 - computation of, 232–238
 - of a partitioned matrix, 298–299
 - reflexive, 243
- generalized quadratic form, 485
- gradient, 275
- Gram–Schmidt orthonormalization, 53, 59–61
- growth curve model, 72–73, 326
- Hölder’s inequality, 446–449
- Hadamard inequality, 333–335, 385
- Hadamard matrix, 369–371
 - normalized, 370
- Hadamard product, 329–339
 - as a Kronecker product, 330
 - eigenvalues of, 336–338
 - nonnegative definite, 332
 - positive definite, 332–333
 - rank of, 330
- Helmert matrix, 16, 33
- Hermite form, 234
- Hermitian matrix, 19
- Hessian, 390
- homogeneous system of equations, 258–260
- hyperplane, 81
- idempotent matrix, 3, 64, 457–462
 - eigenvalues, 457
 - product of, 460
 - rank of, 457
 - sum of, 460
 - symmetric, 459, 461–462
 - trace of, 457
- identity matrix, 2
- indefinite matrix, 17
- independence (linear), 42–45
- independence (stochastic)
 - of quadratic forms, 471–477
 - of random variables, 23
- inequality
 - arithmetic-geometric mean, 452–453
 - Cauchy–Schwarz, 38, 444–446
 - determinantal, 181, 291, 336, 445–447, 450–452, 455–456
 - Hölder’s, 446–449
 - Hadamard, 333–335, 385
 - Jensen’s, 415–417
 - Kantorovich, 455
 - Minkowski’s, 450, 451
 - trace, 376, 440, 445, 447, 449, 451–452, 455–456
 - triangle, 19, 40, 86
- inner product, 38–39
 - Euclidean, 39–40
- interior point, 82
- intersection of vector spaces, 73
- inverse matrix, 9–12
 - and cofactors, 9
 - continuity of, 222
 - derivative of, 398, 401
 - eigenvalues, 100
 - of a partitioned matrix, 285–288
 - of a product, 9
 - of a sum, 10–11
- irreducible matrix, 357

- Jacobi method, 274
- Jacobian matrix, 390–391
- Jensen's inequality, 415–417
- Jordan block matrix, 173
- Jordan decomposition, 173–174
- Kantorovich inequality, 455
- Kronecker product, 315–322
 - determinant of, 319
 - eigenvalues of, 318
 - eigenvectors of, 379
 - generalized inverse of, 318
 - inverse of, 318
 - Moore–Penrose inverse of, 318
 - rank of, 319
 - trace of, 317
- Lagrange function, 417
- Lagrange multipliers, 417
- Lanczos
 - algorithm, 274–278
 - vectors, 276
- latent roots, 96
- latent vectors, 96
- law of cosines, 41
- least squares, 28
 - and best linear unbiased estimator, 130–131
 - and multicollinearity, 109–110, 161–162
 - and solutions to a system of equations, 260–266
 - generalized, 79, 166, 283
 - growth curve model, 72–73, 326
 - in less than full rank models, 64, 266–271
 - in multiple regression, 61–64, 410
 - in multivariate multiple regression, 132, 326
 - in one-way classification model, 90, 269–271, 320–321
 - in ridge regression, 145
 - in simple linear regression, 56–57
 - in two-way classification model, 282, 321–322
 - ordinary, 28–29
 - restricted, 91, 283
 - weighted, 70–71
 - with standardized explanatory variables, 69–70, 109
- least squares inverse, 231
 - computation of, 238
 - least squares solution, 260
- limit point, 81
- linear combination, 36
- linear dependence, 42–45
- linear equations, 71–72
 - and singular value decomposition, 271–273
 - common solution, 281
 - consistency of, 247–251
 - homogeneous system of, 258–260
 - least squares solution of, 260–266
 - linearly independent solutions to, 255–257
 - minimal solution, 280
 - restricted solution, 280
 - solutions to, 251–258
 - sparse systems of, 273–278
 - direct methods, 273
 - iterative methods, 273
 - unique solution to, 254
- linear independence, 42–45
- linear model, 28
- linear space, 36
- linear transformation, 65–73
 - of matrices, 72
- LU factorization, 198
- Mahalanobis distance, 40, 68, 166
- majorization, 433–443
 - and diagonal elements, 437
 - and doubly stochastic matrix, 434–436
 - and eigenvalues, 437–442
 - and order-preserving functions, 442
 - definition, 433–434
- Markov chain, 361–363
- matrix
 - addition, 2
 - backward shift, 367
 - block diagonal, 13
 - circulant, 363–367
 - commutation, 339–346
 - commuting, 182, 184
 - completely positive, 383
 - complex, 18–19
 - correlation, 24, 113–114
 - covariance, 23
 - definition, 1
 - diagonal, 2
 - diagonalizable, 103
 - doubly stochastic, 434
 - duplication, 346–349
 - eigenprojection, 110–112
 - elimination, 349–351
 - forward shift, 367
 - Fourier, 366
 - Hadamard, 369–371
 - Helmert, 16, 33
 - Hermitian, 19
 - Hessian, 390
 - idempotent, 3, 64, 457–461
 - identity, 2
 - indefinite, 17
 - inverse, 9–12
 - irreducible, 357

- Jacobian, 390–391
- Jordan block, 173
- multiplication by a matrix, 3
- multiplication by a scalar, 3
- negative definite, 17
- negative semidefinite, 17
- nilpotent, 149, 195
- nonnegative, 351–363
- nonnegative definite, 17
- nonsingular, 9
- null, 2
- order of, 2
- orthogonal, 15–16
- partitioned, 12–14
- permutation, 16
- positive, 351–357
- positive definite, 17
- positive semidefinite, 17
- primitive, 361
- projection, 58–65, 77–78
- rectangular, 2
- reducible, 357
- semiorthogonal, 16
- similar, 169
- singular, 9
- skew-symmetric, 4
- square, 2
- square root, 17, 163
- symmetric, 4
- Toeplitz, 367–369
- transpose, 3
- triangular, 2
- unitary, 19, 175
- Vandermonde, 371–372
- matrix function, 391
- matrix norm, 184–191
 - Euclidean, 186
 - induced, 185
 - maximum column sum, 186
 - maximum row sum, 186
 - spectral, 186
- maximum
 - absolute, 409
 - conditions for local maximum, 409
 - local, 409
 - of a convex function, 414
 - with equality constraints, 417–423
- maximum likelihood estimation, 411–413, 429
- mean, 20
 - sample, 25
- mean squared error, 192
- mean vector, 23
 - differences in, 118, 134, 178
 - sample, 25
- Minkowski's inequality, 450–451
- minor, 6, 14
 - leading principal, 292, 310
- moment generating function, 21
- moments, 20
- Moore–Penrose inverse, 202
 - and projection matrices, 204, 209
 - and quadratic forms in normal random vectors, 210
 - and the singular value decomposition, 203
 - and the spectral decomposition, 208
 - computation of, 206, 232–234
 - continuity of, 222–224
 - derivative of, 398–400
 - existence of, 202–203
 - of a block diagonal matrix, 219
 - of a diagonal matrix, 208
 - of a matrix product, 211–215
 - of a partitioned matrix, 215–219, 299–302
 - of a sum, 219–221
 - of a symmetric matrix, 207–208
 - properties, 205–211
 - rank, 206
 - uniqueness of, 202–203
- multicollinearity, 109–110, 145, 161–162, 192
- multinomial distribution, 241, 431
- multiplicity of an eigenvalue, 98
 - algebraic, 98
 - geometric, 98
- multivariate normal distribution, 26, 396, 411, 429
 - conditional distribution, 291–292, 295–296
 - density function, 26, 395
 - fourth-order moment matrix, 478, 481
 - maximum likelihood estimates, 411–413
 - moments, 478–480
 - singular, 26
 - sixth-order moment matrix, 500
 - standard, 26
- multivariate t distribution, 27, 34
- negative definite matrix, 17
- negative semidefinite matrix, 17
- nilpotent matrix, 150, 195
- nonnegative definite matrix, 17
 - correlation matrix, 25
 - covariance matrix, 24
- nonnegative matrix, 351–363
 - eigenvalues of, 359–361
 - eigenvectors of, 359–360
 - irreducible, 357
 - primitive, 361
 - reducible, 357
 - spectral radius of, 352
- nonsingular matrix, 9

- norm
 - matrix, 184–191
 - vector, 39, 42
- normal distribution, 21
 - standard, 21
- normalized vector, 15
- null matrix, 2
- null space, 66–67
- null vector, 2
- oblique projection, 76–80
- one-way classification model
 - multivariate, 137–139, 178, 493
 - univariate, 90, 267, 269–271, 320–321, 472, 482
- order
 - of a minor, 14
 - of a square matrix, 2
- order-preserving function, 442
- orthogonal complement, 57–58
 - and null space, 66
 - dimension of, 58
- orthogonal matrix, 15–16
- orthogonal vectors, 15
- orthonormal basis, 53–58
- orthonormal vectors, 15
- parallelogram identity, 86
- partitioned matrix, 12–13
 - determinant of, 288–296
 - eigenvalues of, 302–307
 - generalized inverse of, 298–299
 - inverse of, 285–288
 - Moore–Penrose inverse of, 215–219, 299–302
 - product of, 12
 - rank of, 48, 296–298
- Pearson's chi-squared statistic, 499
- permutation matrix, 16
- perturbation method, 402
 - eigenprojection, 407–408
 - eigenvalue, 403–406, 427
 - matrix inverse, 402–403
 - Moore–Penrose inverse, 426
 - sample correlation matrix, 426
 - symmetric square root, 426
- Poincaré separation theorem, 123
- polar decomposition, 192
- positive definite matrix, 17
 - eigenvalues, 129
 - leading principal minors, 292
- positive matrix, 351–357
 - eigenvalues, 353–357
 - eigenvectors, 353–356
 - spectral radius, 352
- positive semidefinite matrix, 17
 - eigenvalues, 129
- primitive matrix, 361
- principal components analysis, 119–120, 491
- principal submatrix, 124, 292, 310
 - leading, 124
- probability function, 20
 - multivariate, 22
- projection, 53–58
 - oblique, 76–80
 - relative to the A inner product, 79
- projection matrix, 58–65, 209, 231
 - oblique, 77–78
- QR factorization, 60, 165
- quadratic form, 16–17
 - and Moore–Penrose inverse, 210
 - covariance of, 477, 480–481
 - distribution of, 465–471
 - expected value of, 477–485
 - generalized, 485
 - independence of, 471–477
 - matrix of, 17
 - moment generating function of, 498
 - variance of, 477, 480–481
- random variable, 20
 - correlation, 24
 - covariance, 23
 - density function, 20
 - expected value, 20
 - independent, 23
 - mean, 20
 - moment generating function, 21
 - moments, 20
 - probability function, 20
 - variance, 20
- random vector, 22
 - correlation matrix of, 24
 - covariance matrix of, 23
 - density function, 22
 - expected value, 23
 - mean, 23
 - probability function, 22
- range, 45
- rank, 14–15
 - and dimension of null space, 66
 - and eigenvalues, 104, 112, 171–173, 178
 - and linear independence, 45–49
 - full, 14
 - full column, 14
 - full row, 14
 - of a product, 14, 46, 48–49

- of a sum, 46
 - of partitioned matrix, 48, 296–298
- Rayleigh quotient, 117, 120, 275
- reducible matrix, 357
- regression, 28–29
 - best linear unbiased estimator, 130–131
 - best quadratic unbiased estimator, 422–423
 - complete and reduced models, 62, 287, 474
 - F test, 474–475
 - generalized least squares, 79, 166–167, 283
 - minimum variance unbiased linear estimator, 432
 - multicollinearity, 109, 161, 192
 - multiple, 61–64
 - multivariate multiple, 132, 326
 - polynomial, 371
 - principal components, 109–110, 161–162
 - ridge, 145
 - simple linear, 56–57
 - weighted least squares, 70–71
 - with standardized explanatory variables, 69–70, 109–110
- row space, 45
- saddle point, 409
- sample correlation matrix, 25
 - asymptotic covariance matrix of, 495–496
- sample covariance matrix, 25
 - distribution of, 490
 - independent of the sample mean vector, 490
- sample mean, 25
- sample mean vector, 25
- sample variance, 25
 - distribution of, 470–471
 - independent of the sample mean, 475
- Schur complement, 287, 292, 294–295, 306
 - of a Wishart matrix, 487
- Schur decomposition, 175–178
- semiorthogonal matrix, 16
- separating hyperplane theorem, 84
- similar matrices, 169
- simultaneous confidence intervals, 139
- simultaneous diagonalization, 136, 178–184
- singular matrix, 9
- singular value decomposition, 155–162
 - and systems of equations, 271–273
 - of a vector, 159
- singular values, 157
 - and eigenvalues, 160
- skew-symmetric matrix, 4
- spanning set, 37
- spectral decomposition, 108, 111, 162–168
 - of a diagonalizable matrix, 170
- spectral radius, 187
 - of a nonnegative matrix, 352
- spectral set, 111
- spherical distribution, 27
- square root of a matrix, 17, 163
- stationary point, 409
- submatrix, 12–14
- subspace, 36
- sum of squares
 - for error, 29
 - for treatment, 91
- sum of vector spaces, 74
- supporting hyperplane theorem, 83
- symmetric matrix, 4
- T -transform, 435, 437
- Taylor formula
 - first-order, 388–389
 - for a vector function, 390
 - k th-order, 388–389
- time series, 368
- Toeplitz matrix, 367–369
- trace, 4
 - and determinant, 452
 - and eigenvalues, 101
 - derivative of, 396
 - of a product, 4, 440, 445, 448, 449
 - of a sum, 451
- trace inequality, 376, 440, 445, 447, 449, 451, 452, 455, 456
- transition probabilities, 361
- transpose, 3–4
 - conjugate, 19
- transpose product, 13, 167
 - eigenvalues, 131
- triangle inequality, 19, 40, 86
- triangular matrix
 - definition, 2
 - lower, 2
 - upper, 2
- two-way classification model, 282, 321–322, 374
- uniform distribution, 27
 - fourth-order moment matrix, 501
- union-intersection procedure, 138–139, 431
- unit vector, 15
- unitary matrix, 19, 175
- univariate normal distribution, 21
 - standard, 21
- Vandermonde matrix, 371–372
- variance, 20
 - of a quadratic form, 477, 480–481
 - sample, 25
- vec operator, 324–328

- vector, 2
 - column, 2
 - definition, 2
 - normalized, 15
 - null, 2
 - orthogonal, 15
 - orthonormal, 15
 - row, 2
 - unit, 15
- vector norm, 39
 - Euclidean, 40, 42
 - infinity norm, 42
 - max norm, 42
 - sum norm, 42
- vector space, 36
 - basis of, 49–52
 - definition, 35
 - dimension of, 49
 - direct sum, 75
 - Euclidean, 40
 - intersection, 73
 - projection matrix of, 59–65, 77–78
 - spanning set, 37
 - sum, 74
- vector subspace, 36
- Weyl's Theorem, 124
- Wishart distribution, 485–490
 - and sample covariance matrix, 490–491
 - covariance matrix of, 488
 - mean of, 488

WILEY SERIES IN PROBABILITY AND STATISTICS

established by Walter A. Shewhart and Samuel S. Wilks

Editors: David J. Balding, Noel A. C. Cressie, Garrett M. Fitzmaurice,
Geof H. Givens, Harvey Goldstein, Geert Molenberghs, David W. Scott,
Adrian F. M. Smith, Ruey S. Tsay, Sanford Weisberg
Editors Emeriti: J. Stuart Hunter, Iain M. Johnstone, Joseph B. Kadane,
Jozef L. Teugels

The Wiley Series in Probability and Statistics is well established and authoritative. It covers many topics of current research interest in both pure and applied statistics and probability theory. Written by leading statisticians and institutions, the titles span both state-of-the-art developments in the field and classical methods.

Reflecting the wide range of current research in statistics, the series encompasses applied, methodological and theoretical statistics, ranging from applications and new techniques made possible by advances in computerized practice to rigorous treatment of theoretical approaches.

This series provides essential and invaluable reading for all statisticians, whether in academia, industry, government, or research.

† ABRAHAM and LEDOLTER · Statistical Methods for Forecasting

AGRESTI · Analysis of Ordinal Categorical Data, Second Edition

AGRESTI · An Introduction to Categorical Data Analysis, Second Edition

AGRESTI · Categorical Data Analysis, Third Edition

AGRESTI · Foundations of Linear and Generalized Linear Models

ALSTON, Mengersen and Pettitt (editors) · Case Studies in Bayesian Statistical Modelling and Analysis

ALTMAN, GILL, and McDONALD · Numerical Issues in Statistical Computing for the Social Scientist

AMARATUNGA and CABRERA · Exploration and Analysis of DNA Microarray and Protein Array Data

AMARATUNGA, CABRERA, and SHKEDY · Exploration and Analysis of DNA Microarray and Other High-Dimensional Data, Second Edition

ANDÉL · Mathematics of Chance

ANDERSON · An Introduction to Multivariate Statistical Analysis, Third Edition

* ANDERSON · The Statistical Analysis of Time Series

ANDERSON, AUQUIER, HAUCK, OAKES, VANDAELE, and WEISBERG · Statistical Methods for Comparative Studies

ANDERSON and LOYNES · The Teaching of Practical Statistics

ARMITAGE and DAVID (editors) · Advances in Biometry

ARNOLD, BALAKRISHNAN, and NAGARAJA · Records

* ARTHANARI and DODGE · Mathematical Programming in Statistics

AUGUSTIN, COOLEN, DE COOMAN and TROFFAES (editors) · Introduction to Imprecise Probabilities

* BAILEY · The Elements of Stochastic Processes with Applications to the Natural Sciences

BAJORSKI · Statistics for Imaging, Optics, and Photonics

BALAKRISHNAN and KOUTRAS · Runs and Scans with Applications

BALAKRISHNAN and NG · Precedence-Type Tests and Applications

BALI, ENGLE, and MURRAY · Empirical Asset Pricing: The Cross-Section of Stock Returns

BARNETT · Comparative Statistical Inference, Third Edition

BARNETT · Environmental Statistics

BARNETT and LEWIS · Outliers in Statistical Data, Third Edition

Bartholomew, Knott, and Moustaki · Latent Variable Models and Factor Analysis: A Unified Approach, Third Edition

BARTOSZYNSKI and NIEWIADOMSKA-BUGAJ · Probability and Statistical Inference, Second Edition

BASILEVSKY · Statistical Factor Analysis and Related Methods: Theory and Applications

BATES and WATTS · Nonlinear Regression Analysis and Its Applications

BECHHOFFER, SANTNER, and GOLDSMAN · Design and Analysis of Experiments for Statistical Selection, Screening, and Multiple Comparisons

BEH and LOMBARDO · Correspondence Analysis: Theory, Practice and New Strategies

BEIRLANT, GOEGBEUR, SEGERS, TEUGELS, and DE WAAL · Statistics of Extremes: Theory and Applications

BELSLEY · Conditioning Diagnostics: Collinearity and Weak Data in Regression

† BELSLEY, KUH, and WELSCH · Regression Diagnostics: Identifying Influential Data and Sources of Collinearity

BENDAT and PIERSON · Random Data: Analysis and Measurement Procedures, Fourth Edition

BERNARDO and SMITH · Bayesian Theory

BHAT and MILLER · Elements of Applied Stochastic Processes, Third Edition

BHATTACHARYA and WAYMIRE · Stochastic Processes with Applications

BIEMER, GROVES, LYBERG, MATHIOWETZ, and SUDMAN · Measurement Errors in Surveys

BILLINGSLEY · Convergence of Probability Measures, Second Edition

BILLINGSLEY · Probability and Measure, Anniversary Edition

BIRKES and DODGE · Alternative Methods of Regression

Bisgaard and Kulahci · Time Series Analysis and Forecasting by Example

Biswas, Datta, Fine, and Segal · Statistical Advances in the Biomedical Sciences: Clinical Trials, Epidemiology, Survival Analysis, and Bioinformatics

BLISCHKE and MURTHY (editors) · Case Studies in Reliability and Maintenance

BLISCHKE and MURTHY · Reliability: Modeling, Prediction, and Optimization

BLOOMFIELD · Fourier Analysis of Time Series: An Introduction, Second Edition

BOLLEN · Structural Equations with Latent Variables

BOLLEN and CURRAN · Latent Curve Models: A Structural Equation Perspective

BONNINI, CORAIN, MAROZZI and SALMASO · Nonparametric Hypothesis Testing: Rank and Permutation Methods with Applications in R

BOROVKOV · Ergodicity and Stability of Stochastic Processes

BOSQ and BLANKE · Inference and Prediction in Large Dimensions

BOULEAU · Numerical Methods for Stochastic Processes

* BOX and TIAO · Bayesian Inference in Statistical Analysis

BOX · Improving Almost Anything, Revised Edition

* BOX and DRAPER · Evolutionary Operation: A Statistical Method for Process Improvement

BOX and DRAPER · Response Surfaces, Mixtures, and Ridge Analyses, Second Edition

BOX, HUNTER, and HUNTER · Statistics for Experimenters: Design, Innovation, and Discovery, Second Edition

BOX, JENKINS, REINSEL, and LJUNG · Time Series Analysis: Forecasting and Control, Fifth Edition

BOX, LUCEÑO, and Paniagua-QuiÑones · Statistical Control by Monitoring and Adjustment, Second Edition

* BROWN and HOLLANDER · Statistics: A Biomedical Introduction

CAIROLI and DALANG · Sequential Stochastic Optimization

CASTILLO, HADI, BALAKRISHNAN, and SARABIA · Extreme Value and Related Models with Applications in Engineering and Science

CHAN · Time Series: Applications to Finance with R and S-Plus®, Second Edition

CHARALAMBIDES · Combinatorial Methods in Discrete Distributions

CHATTERJEE and HADI · Regression Analysis by Example, Fourth Edition

CHATTERJEE and HADI · Sensitivity Analysis in Linear Regression

Chen · The Fitness of Information: Quantitative Assessments of Critical Evidence

CHERNICK · Bootstrap Methods: A Guide for Practitioners and Researchers, Second Edition

CHERNICK and FRIIS · Introductory Biostatistics for the Health Sciences

CHILÈS and DELFINER · Geostatistics: Modeling Spatial Uncertainty, Second Edition

CHIU, STOYAN, KENDALL and MECKE · Stochastic Geometry and Its Applications, Third Edition

CHOW and LIU · Design and Analysis of Clinical Trials: Concepts and Methodologies, Third Edition

CLARKE · Linear Models: The Theory and Application of Analysis of Variance

CLARKE and DISNEY · Probability and Random Processes: A First Course with Applications, Second Edition

* COCHRAN and COX · Experimental Designs, Second Edition

COLLINS and LANZA · Latent Class and Latent Transition Analysis: With Applications in the Social, Behavioral, and Health Sciences

CONGDON · Applied Bayesian Modelling, Second Edition

CONGDON · Bayesian Models for Categorical Data

CONGDON · Bayesian Statistical Modelling, Second Edition

CONOVER · Practical Nonparametric Statistics, Third Edition

COOK · Regression Graphics

COOK and WEISBERG · An Introduction to Regression Graphics

COOK and WEISBERG · Applied Regression Including Computing and Graphics

CORNELL · A Primer on Experiments with Mixtures

CORNELL · Experiments with Mixtures, Designs, Models, and the Analysis of Mixture Data, Third Edition

COX · A Handbook of Introductory Statistical Methods

CRESSIE · Statistics for Spatial Data, Revised Edition

CRESSIE and WIKLE · Statistics for Spatio-Temporal Data

CSÖRGÖ and HORVÁTH · Limit Theorems in Change Point Analysis

Dagpunar · Simulation and Monte Carlo: With Applications in Finance and MCMC

DANIEL · Applications of Statistics to Industrial Experimentation

DANIEL · Biostatistics: A Foundation for Analysis in the Health Sciences, Eighth Edition

* DANIEL · Fitting Equations to Data: Computer Analysis of Multifactor Data, Second Edition

DASU and JOHNSON · Exploratory Data Mining and Data Cleaning

DAVID and NAGARAJA · Order Statistics, Third Edition

DAVINO, FURNO and VISTOCCO · Quantile Regression: Theory and Applications

* DEGROOT, FIENBERG, and KADANE · Statistics and the Law

DEL CASTILLO · Statistical Process Adjustment for Quality Control

DeMaris · Regression with Social Data: Modeling Continuous and Limited Response Variables

DEMIDENKO · Mixed Models: Theory and Applications with R, Second Edition

Denison, Holmes, Mallick, and Smith · Bayesian Methods for Nonlinear Classification and Regression

DETTE and STUDDEN · The Theory of Canonical Moments with Applications in Statistics, Probability, and Analysis

DEY and MUKERJEE · Fractional Factorial Plans

DILLON and GOLDSTEIN · Multivariate Analysis: Methods and Applications

* DODGE and ROMIG · Sampling Inspection Tables, Second Edition

* DOOB · Stochastic Processes

DOWDY, WEARDEN, and CHILKO · Statistics for Research, Third Edition

DRAPER and SMITH · Applied Regression Analysis, Third Edition

DRYDEN and MARDIA · Statistical Shape Analysis

DUDEWICZ and MISHRA · Modern Mathematical Statistics

DUNN and CLARK · Basic Statistics: A Primer for the Biomedical Sciences, Fourth Edition

DUPUIS and ELLIS · A Weak Convergence Approach to the Theory of Large Deviations

EDLER and KITSOS · Recent Advances in Quantitative Methods in Cancer and Human Health Risk Assessment

* ELANDT-JOHNSON and JOHNSON · Survival Models and Data Analysis

ENDERS · Applied Econometric Time Series, Third Edition

† ETHIER and KURTZ · Markov Processes: Characterization and Convergence

EVANS, HASTINGS, and PEACOCK · Statistical Distributions, Third Edition

EVERITT, LANDAU, LEESE, and STAHL · Cluster Analysis, Fifth Edition

FEDERER and KING · Variations on Split Plot and Split Block Experiment Designs

FELLER · An Introduction to Probability Theory and Its Applications, Volume I, Third Edition, Revised; Volume II, Second Edition

FITZMAURICE, LAIRD, and WARE · Applied Longitudinal Analysis, Second Edition

* FLEISS · The Design and Analysis of Clinical Experiments

FLEISS · Statistical Methods for Rates and Proportions, Third Edition

† FLEMING and HARRINGTON · Counting Processes and Survival Analysis

FUJIKOSHI, ULYANOV, and SHIMIZU · Multivariate Statistics: High-Dimensional and Large-Sample Approximations

FULLER · Introduction to Statistical Time Series, Second Edition

† FULLER · Measurement Error Models

GALLANT · Nonlinear Statistical Models

GEISSER · Modes of Parametric Statistical Inference

GELMAN and MENG · Applied Bayesian Modeling and Causal Inference from ncomplete-Data Perspectives

GEWEKE · Contemporary Bayesian Econometrics and Statistics

GHOSH, MUKHOPADHYAY, and SEN · Sequential Estimation

GIESBRECHT and GUMPERTZ · Planning, Construction, and Statistical Analysis of Comparative Experiments

GIFI · Nonlinear Multivariate Analysis

GIVENS and HOETING · Computational Statistics

GLASSERMAN and YAO · Monotone Structure in Discrete-Event Systems

GNANADESIKAN · Methods for Statistical Data Analysis of Multivariate Observations, Second Edition

GOLDSTEIN · Multilevel Statistical Models, Fourth Edition

GOLDSTEIN and LEWIS · Assessment: Problems, Development, and Statistical Issues

Goldstein and Wooff · Bayes Linear Statistics

GRAHAM · Markov Chains: Analytic and Monte Carlo Computations

GREENWOOD and NIKULIN · A Guide to Chi-Squared Testing

GROSS, SHORTLE, THOMPSON, and HARRIS · Fundamentals of Queueing Theory, Fourth Edition

GROSS, SHORTLE, THOMPSON, and HARRIS · Solutions Manual to Accompany Fundamentals of Queueing Theory, Fourth Edition

* HAHN and SHAPIRO · Statistical Models in Engineering

HAHN and MEEKER · Statistical Intervals: A Guide for Practitioners

HALD · A History of Probability and Statistics and their Applications Before 1750

† HAMPEL · Robust Statistics: The Approach Based on Influence Functions

Hartung, Knapp, and Sinha · Statistical Meta-Analysis with Applications

HEIBERGER · Computation for the Analysis of Designed Experiments

HEDAYAT and SINHA · Design and Inference in Finite Population Sampling

HEDEKER and GIBBONS · Longitudinal Data Analysis

HELLER · MACSYMA for Statisticians

HERITIER, CANTONI, COPT, and VICTORIA-FESER · Robust Methods in Biostatistics

HINKELMANN and KEMPTHORNE · Design and Analysis of Experiments, Volume 1: Introduction to Experimental Design, Second Edition

HINKELMANN and KEMPTHORNE · Design and Analysis of Experiments, Volume 2: Advanced Experimental Design

HINKELMANN (editor) · Design and Analysis of Experiments, Volume 3: Special Designs and Applications

HOAGLIN, MOSTELLER, and TUKEY · Fundamentals of Exploratory Analysis of Variance

* HOAGLIN, MOSTELLER, and TUKEY · Exploring Data Tables, Trends and Shapes

* HOAGLIN, MOSTELLER, and TUKEY · Understanding Robust and Exploratory Data Analysis

HOCHBERG and TAMHANE · Multiple Comparison Procedures

HOCKING · Methods and Applications of Linear Models: Regression and the Analysis of Variance, Third Edition

HOEL · Introduction to Mathematical Statistics, Fifth Edition

HOGG and KLUGMAN · Loss Distributions

HOLLANDER, WOLFE, and CHICKEN · Nonparametric Statistical Methods, Third Edition

HOSMER and LEMESHOW · Applied Logistic Regression, Second Edition

HOSMER, LEMESHOW, and MAY · Applied Survival Analysis: Regression Modeling of Time-to-Event Data, Second Edition

HUBER · Data Analysis: What Can Be Learned From the Past 50 Years

HUBER · Robust Statistics

† HUBER and Ronchetti · Robust Statistics, Second Edition

HUBERTY · Applied Discriminant Analysis, Second Edition

HUBERTY and OLEJNIK · Applied MANOVA and Discriminant Analysis, Second Edition

HUITEMA · The Analysis of Covariance and Alternatives: Statistical Methods for Experiments, Quasi-Experiments, and Single-Case Studies, Second Edition

HUNT and KENNEDY · Financial Derivatives in Theory and Practice, Revised Edition

HURD and MIAMEE · Periodically Correlated Random Sequences: Spectral Theory and Practice

HUSKOVA, BERAN, and DUPAC · Collected Works of Jaroslav Hajek—with Commentary

HUZURBAZAR · Flowgraph Models for Multistate Time-to-Event Data

Jackman · Bayesian Analysis for the Social Sciences

† JACKSON · A User's Guide to Principle Components

JOHN · Statistical Methods in Engineering and Quality Assurance

JOHNSON · Multivariate Statistical Simulation

JOHNSON and BALAKRISHNAN · Advances in the Theory and Practice of Statistics: A Volume in Honor of Samuel Kotz

JOHNSON, KEMP, and KOTZ · Univariate Discrete Distributions, Third Edition

JOHNSON and KOTZ (editors) · Leading Personalities in Statistical Sciences: From the Seventeenth Century to the Present

JOHNSON, KOTZ, and BALAKRISHNAN · Continuous Univariate Distributions, Volume 1, Second Edition

JOHNSON, KOTZ, and BALAKRISHNAN · Continuous Univariate Distributions, Volume 2, Second Edition

JOHNSON, KOTZ, and BALAKRISHNAN · Discrete Multivariate Distributions

JUDGE, GRIFFITHS, HILL, LÜTKEPOHL, and LEE · The Theory and Practice of Econometrics, Second Edition

JUREK and MASON · Operator-Limit Distributions in Probability Theory

KADANE · Bayesian Methods and Ethics in a Clinical Trial Design

KADANE AND SCHUM · A Probabilistic Analysis of the Sacco and Vanzetti Evidence

KALBFLEISCH and PRENTICE · The Statistical Analysis of Failure Time Data, Second Edition

KARIYA and KURATA · Generalized Least Squares

KASS and VOS · Geometrical Foundations of Asymptotic Inference

† KAUFMAN and ROUSSEUW · Finding Groups in Data: An Introduction to Cluster Analysis

KEDEM and FOKIANOS · Regression Models for Time Series Analysis

KENDALL, BARDEN, CARNE, and LE · Shape and Shape Theory

KHURI · Advanced Calculus with Applications in Statistics, Second Edition

KHURI, MATHEW, and SINHA · Statistical Tests for Mixed Linear Models

* KISH · Statistical Design for Research

KLEIBER and KOTZ · Statistical Size Distributions in Economics and Actuarial Sciences

Klemelä · Smoothing of Multivariate Data: Density Estimation and Visualization
 KLUGMAN, PANJER, and WILLMOT · Loss Models: From Data to Decisions, Third Edition
 KLUGMAN, PANJER, and WILLMOT · Loss Models: Further Topics
 KLUGMAN, PANJER, and WILLMOT · Solutions Manual to Accompany Loss Models: From Data to Decisions, Third Edition
 KOSKI and NOBLE · Bayesian Networks: An Introduction
 KOTZ, BALAKRISHNAN, and JOHNSON · Continuous Multivariate Distributions, Volume 1, Second Edition
 KOTZ and JOHNSON (editors) · Encyclopedia of Statistical Sciences: Volumes 1 to 9 with Index
 KOTZ and JOHNSON (editors) · Encyclopedia of Statistical Sciences: Supplement Volume
 KOTZ, READ, and BANKS (editors) · Encyclopedia of Statistical Sciences: Update Volume 1
 KOTZ, READ, and BANKS (editors) · Encyclopedia of Statistical Sciences: Update Volume 2
 KOWALSKI and TU · Modern Applied U-Statistics
 Krishnamoorthy and Mathew · Statistical Tolerance Regions: Theory, Applications, and Computation
 Kroese, Taimre, and Botev · Handbook of Monte Carlo Methods
 KROONENBERG · Applied Multiway Data Analysis
 KULINSKAYA, MORGENTHALER, and STAUDTE · Meta Analysis: A Guide to Calibrating and Combining Statistical Evidence
 Kulkarni and Harman · An Elementary Introduction to Statistical Learning Theory
 KUROWICKA and COOKE · Uncertainty Analysis with High Dimensional Dependence Modelling
 KVAM and VIDAKOVIC · Nonparametric Statistics with Applications to Science and Engineering
 LACHIN · Biostatistical Methods: The Assessment of Relative Risks, Second Edition
 LAD · Operational Subjective Statistical Methods: A Mathematical, Philosophical, and Historical Introduction
 LAMPERTI · Probability: A Survey of the Mathematical Theory, Second Edition
 LAWLESS · Statistical Models and Methods for Lifetime Data, Second Edition
 LAWSON · Statistical Methods in Spatial Epidemiology, Second Edition
 LE · Applied Categorical Data Analysis, Second Edition
 LE · Applied Survival Analysis
 Lee · Structural Equation Modeling: A Bayesian Approach
 LEE and WANG · Statistical Methods for Survival Data Analysis, Fourth Edition
 LePAGE and BILLARD · Exploring the Limits of Bootstrap
 LESSLER and KALSBECK · Nonsampling Errors in Surveys
 LEYLAND and GOLDSTEIN (editors) · Multilevel Modelling of Health Statistics
 LIAO · Statistical Group Comparison
 LIN · Introductory Stochastic Analysis for Finance and Insurance
 LINDLEY · Understanding Uncertainty, Revised Edition
 LITTLE and RUBIN · Statistical Analysis with Missing Data, Second Edition
 Lloyd · The Statistical Analysis of Categorical Data
 LOWEN and TEICH · Fractal-Based Point Processes
 MAGNUS and NEUDECKER · Matrix Differential Calculus with Applications in Statistics and Econometrics, Revised Edition
 MALLER and ZHOU · Survival Analysis with Long Term Survivors

MARCHETTE · Random Graphs for Statistical Pattern Recognition

MARDIA and JUPP · Directional Statistics

MARKOVICH · Nonparametric Analysis of Univariate Heavy-Tailed Data: Research and Practice

MARONNA, MARTIN and YOHAI · Robust Statistics: Theory and Methods

MASON, GUNST, and HESS · Statistical Design and Analysis of Experiments with Applications to Engineering and Science, Second Edition

McCULLOCH, SEARLE, and NEUHAUS · Generalized, Linear, and Mixed Models, Second Edition

McFADDEN · Management of Data in Clinical Trials, Second Edition

* McLACHLAN · Discriminant Analysis and Statistical Pattern Recognition

McLACHLAN, DO, and AMBROISE · Analyzing Microarray Gene Expression Data

McLACHLAN and KRISHNAN · The EM Algorithm and Extensions, Second Edition

McLACHLAN and PEEL · Finite Mixture Models

McNEIL · Epidemiological Research Methods

Meeker and Escobar · Statistical Methods for Reliability Data

MEERSCHAERT and SCHEFFLER · Limit Distributions for Sums of Independent Random Vectors: Heavy Tails in Theory and Practice

Mengersen, Robert, and Titterton · Mixtures: Estimation and Applications

MICKEY, DUNN, and CLARK · Applied Statistics: Analysis of Variance and Regression, Third Edition

* MILLER · Survival Analysis, Second Edition

MONTGOMERY, JENNINGS, and KULAHCI · Introduction to Time Series Analysis and Forecasting, Second Edition

MONTGOMERY, PECK, and VINING · Introduction to Linear Regression Analysis, Fifth Edition

MORGENTHALER and TUKEY · Configural Polysampling: A Route to Practical Robustness

MUIRHEAD · Aspects of Multivariate Statistical Theory

Muller and Stoyan · Comparison Methods for Stochastic Models and Risks

MURTHY, XIE, and JIANG · Weibull Models

MYERS, MONTGOMERY, and ANDERSON-COOK · Response Surface Methodology: Process and Product Optimization Using Designed Experiments, Third Edition

MYERS, MONTGOMERY, VINING, and ROBINSON · Generalized Linear Models. With Applications in Engineering and the Sciences, Second Edition

Natvig · Multistate Systems Reliability Theory With Applications

† NELSON · Accelerated Testing, Statistical Models, Test Plans, and Data Analyses

† NELSON · Applied Life Data Analysis

NEWMAN · Biostatistical Methods in Epidemiology

Ng, Tain, and Tang · Dirichlet Theory: Theory, Methods and Applications

OKABE, BOOTS, SUGIHARA, and CHIU · Spatial Tessellations: Concepts and Applications of Voronoi Diagrams, Second Edition

OLIVER and SMITH · Influence Diagrams, Belief Nets and Decision Analysis

PALTA · Quantitative Methods in Population Health: Extensions of Ordinary Regressions

PANJER · Operational Risk: Modeling and Analytics

PANKRATZ · Forecasting with Dynamic Regression Models

PANKRATZ · Forecasting with Univariate Box-Jenkins Models: Concepts and Cases

PARDOUX · Markov Processes and Applications: Algorithms, Networks, Genome and Finance

Parmigiani and Inoue · Decision Theory: Principles and Approaches

* PARZEN · Modern Probability Theory and Its Applications

PEÑA, TIAO, and TSAY · A Course in Time Series Analysis

Pesarin and Salmaso · Permutation Tests for Complex Data: Applications and Software

PIANTADOSI · Clinical Trials: A Methodologic Perspective, Second Edition

POURAHMADI · Foundations of Time Series Analysis and Prediction Theory

POURAHMADI · High-Dimensional Covariance Estimation

POWELL · Approximate Dynamic Programming: Solving the Curses of Dimensionality, Second Edition

POWELL and RYZHOV · Optimal Learning

PRESS · Subjective and Objective Bayesian Statistics, Second Edition

PRESS and TANUR · The Subjectivity of Scientists and the Bayesian Approach

PURI, VILAPLANA, and WERTZ · New Perspectives in Theoretical and Applied Statistics

† PUTERMAN · Markov Decision Processes: Discrete Stochastic Dynamic Programming

QIU · Image Processing and Jump Regression Analysis

* RAO · Linear Statistical Inference and Its Applications, Second Edition

RAO · Statistical Inference for Fractional Diffusion Processes

RAUSAND and HØYLAND · System Reliability Theory: Models, Statistical Methods, and Applications, Second Edition

Rayner, THAS, and BEST · Smooth Tests of Goodness of Fit: Using R, Second Edition

RENCHEr and SCHAALJE · Linear Models in Statistics, Second Edition

RENCHEr and CHRISTENSEN · Methods of Multivariate Analysis, Third Edition

RENCHEr · Multivariate Statistical Inference with Applications

RIGDON and BASU · Statistical Methods for the Reliability of Repairable Systems

* RIPLEY · Spatial Statistics

* RIPLEY · Stochastic Simulation

ROHATGI and SALEH · An Introduction to Probability and Statistics, Third Edition

ROLSKI, SCHMIDT, SCHMIDT, and TEUGELS · Stochastic Processes for Insurance and Finance

ROSENBERGER and LACHIN · Randomization in Clinical Trials: Theory and Practice

ROSSI, ALLENBY, and McCULLOCH · Bayesian Statistics and Marketing

† ROUSSEUW and LEROY · Robust Regression and Outlier Detection

Royston and Sauerbrei · Multivariate Model Building: A Pragmatic Approach to Regression Analysis Based on Fractional Polynomials for Modeling Continuous Variables

* RUBIN · Multiple Imputation for Nonresponse in Surveys

RUBINSTEIN and KROESE · Simulation and the Monte Carlo Method, Second Edition

RUBINSTEIN and MELAMED · Modern Simulation and Modeling

RUBINSTEIN, RIDDER, and VAISMAN · Fast Sequential Monte Carlo Methods for Counting and Optimization

RYAN · Modern Engineering Statistics

RYAN · Modern Experimental Design

RYAN · Modern Regression Methods, Second Edition

Ryan · Sample Size Determination and Power

RYAN · Statistical Methods for Quality Improvement, Third Edition

SALEH · Theory of Preliminary Test and Stein-Type Estimation with Applications

SALTELLI, CHAN, and SCOTT (editors) · Sensitivity Analysis

Scherer · Batch Effects and Noise in Microarray Experiments: Sources and Solutions

* SCHEFFE · The Analysis of Variance

SCHIMEK · Smoothing and Regression: Approaches, Computation, and Application

SCHOTT · Matrix Analysis for Statistics, Second Edition

SCHOTT · Matrix Analysis for Statistics, Third Edition

Schoutens · Levy Processes in Finance: Pricing Financial Derivatives

SCOTT · Multivariate Density Estimation

SCOTT · Multivariate Density Estimation: Theory, Practice, and Visualization

* SEARLE · Linear Models

† SEARLE · Linear Models for Unbalanced Data

† SEARLE · Matrix Algebra Useful for Statistics

† SEARLE, CASELLA, and McCULLOCH · Variance Components

SEARLE and WILLETT · Matrix Algebra for Applied Economics

SEBER · A Matrix Handbook For Statisticians

† SEBER · Multivariate Observations

SEBER and LEE · Linear Regression Analysis, Second Edition

† SEBER and WILD · Nonlinear Regression

SENNOTT · Stochastic Dynamic Programming and the Control of Queueing Systems

* SERFLING · Approximation Theorems of Mathematical Statistics

SHAFER and VOVK · Probability and Finance: It's Only a Game!

SHERMAN · Spatial Statistics and Spatio-Temporal Data: Covariance Functions and Directional Properties

SILVAPULLE and SEN · Constrained Statistical Inference: Inequality, Order, and Shape Restrictions

SINGPURWALLA · Reliability and Risk: A Bayesian Perspective

SMALL and McLEISH · Hilbert Space Methods in Probability and Statistical Inference

SRIVASTAVA · Methods of Multivariate Statistics

STAPLETON · Linear Statistical Models, Second Edition

STAPLETON · Models for Probability and Statistical Inference: Theory and Applications

STAUDTE and SHEATHER · Robust Estimation and Testing

Stoyan · Counterexamples in Probability, Second Edition

STOYAN and STOYAN · Fractals, Random Shapes and Point Fields: Methods of Geometrical Statistics

STREET and BURGESS · The Construction of Optimal Stated Choice Experiments: Theory and Methods

STYAN · The Collected Papers of T. W. Anderson: 1943–1985

SUTTON, ABRAMS, JONES, SHELDON, and SONG · Methods for Meta-Analysis in Medical Research

TAKEZAWA · Introduction to Nonparametric Regression

TAMHANE · Statistical Analysis of Designed Experiments: Theory and Applications

TANAKA · Time Series Analysis: Nonstationary and Noninvertible Distribution Theory

THOMPSON · Empirical Model Building: Data, Models, and Reality, Second Edition

THOMPSON · Sampling, Third Edition

THOMPSON · Simulation: A Modeler's Approach

THOMPSON and SEBER · Adaptive Sampling

THOMPSON, WILLIAMS, and FINDLAY · Models for Investors in Real World Markets

TIERNEY · LISP-STAT: An Object-Oriented Environment for Statistical Computing and Dynamic Graphics

TROFFAES and DE COOMAN · Lower Previsions

TSAY · Analysis of Financial Time Series, Third Edition

TSAY · An Introduction to Analysis of Financial Data with R

TSAY · Multivariate Time Series Analysis: With R and Financial Applications

UPTON and FINGLETON · Spatial Data Analysis by Example, Volume II: Categorical and Directional Data

† VAN BELLE · Statistical Rules of Thumb, Second Edition

VAN BELLE, FISHER, HEAGERTY, and LUMLEY · Biostatistics: A Methodology for the Health Sciences, Second Edition

VESTRUP · The Theory of Measures and Integration

VIDAKOVIC · Statistical Modeling by Wavelets

Viertl · Statistical Methods for Fuzzy Data

VINOD and REAGLE · Preparing for the Worst: Incorporating Downside Risk in Stock Market Investments

WALLER and GOTWAY · Applied Spatial Statistics for Public Health Data

WEISBERG · Applied Linear Regression, Fourth Edition

WEISBERG · Bias and Causation: Models and Judgment for Valid Comparisons

WELSH · Aspects of Statistical Inference

WESTFALL and YOUNG · Resampling-Based Multiple Testing: Examples and Methods for p-Value Adjustment

* WHITTAKER · Graphical Models in Applied Multivariate Statistics

WINKER · Optimization Heuristics in Economics: Applications of Threshold Accepting

WOODWORTH · Biostatistics: A Bayesian Introduction

WOOLSON and CLARKE · Statistical Methods for the Analysis of Biomedical Data, Second Edition

WU and HAMADA · Experiments: Planning, Analysis, and Parameter Design Optimization, Second Edition

WU and ZHANG · Nonparametric Regression Methods for Longitudinal Data Analysis

Yakir · Extremes in Random Fields

YIN · Clinical Trial Design: Bayesian and Frequentist Adaptive Methods

YOUNG, VALERO-MORA, and FRIENDLY · Visual Statistics: Seeing Data with Dynamic Interactive Graphics

ZACKS · Examples and Problems in Mathematical Statistics

ZACKS · Stage-Wise Adaptive Designs

* ZELLNER · An Introduction to Bayesian Inference in Econometrics

ZELTERMAN · Discrete Distributions—Applications in the Health Sciences

ZHOU, OBUCHOWSKI, and MCCLISH · Statistical Methods in Diagnostic Medicine, Second Edition

WILEY END USER LICENSE AGREEMENT

Go to www.wiley.com/go/eula to access Wiley's ebook EULA.