

---

# Matrix Analysis for Statistics

James R. Schott



---

WILEY SERIES IN  
PROBABILITY AND STATISTICS

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: *Vic Barnett, Ralph A. Bradley, Nicholas I. Fisher, J. Stuart Hunter, J. B. Kadane, David G. Kendall, David W. Scott, Adrian F. M. Smith, Jozef L. Teugels, Geoffrey S. Watson*

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

# Matrix Analysis for Statistics

JAMES R. SCHOTT



A Wiley-Interscience Publication

JOHN WILEY & SONS, INC.

New York • Chichester • Brisbane • Toronto • Singapore • Weinheim

This text is printed on acid-free paper.

Copyright © 1997 By John Wiley & Sons, Inc.

All Rights Reserved. Published Simultaneously In Canada.

Reproduction or translation of any part of this work beyond that permitted by Section 107 or 108 of the 1976 United States Copyright Act without the permission of the copyright owner is unlawful. Requests for permission or further information should be addressed to the Permissions Department, John Wiley & Sons, Inc., 605 Third Avenue, New York, NY 10158-0012.

***Library Of Congress Cataloging In Publication Data:***

Schott, James R., 1955-

Matrix Analysis For Statistics / James R. Schott.

p. cm. – (Wiley Series In Probability And Statistics. Applied Probability And Statistics)

“A Wiley-Interscience Publication.”

Includes bibliographical references and index.

ISBN 0-471-15409-1 (cloth : alk. paper)

1. Matrices. 2. Mathematical Statistics. I. Title. II. Series.

QA188.S24 1996

512.9'434—dc20

96-12133

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

To  
Susan, Adam, and Sarah



# Contents

## Preface

<b>1. A Review of Elementary Matrix Algebra</b>	<b>1</b>
1. Introduction, 1	
2. Definitions and Notation, 1	
3. Matrix Addition and Multiplication, 2	
4. The Transpose, 3	
5. The Trace, 4	
6. The Determinant, 5	
7. The Inverse, 8	
8. Partitioned Matrices, 11	
9. The Rank of a Matrix, 13	
10. Orthogonal Matrices, 14	
11. Quadratic Forms, 15	
12. Complex Matrices, 16	
13. Random Vectors and Some Related Statistical Concepts, 18	
Problems, 28	
<b>2. Vector Spaces</b>	<b>32</b>
1. Introduction, 32	
2. Definitions, 32	
3. Linear Independence and Dependence, 38	
4. Bases and Dimension, 41	
5. Matrix Rank and Linear Independence, 43	
6. Orthonormal Bases and Projections, 48	
7. Projection Matrices, 52	
8. Linear Transformations and Systems of Linear Equations, 60	

9.	The Intersection and Sum of Vector Spaces, 67	
10.	Convex Sets, 70	
	Problems, 74	
<b>3.</b>	<b>Eigenvalues and Eigenvectors</b>	<b>84</b>
1.	Introduction, 84	
2.	Eigenvalues, Eigenvectors, and Eigenspaces, 84	
3.	Some Basic Properties of Eigenvalues and Eigenvectors, 88	
4.	Symmetric Matrices, 93	
5.	Continuity of Eigenvalues and Eigenprojections, 102	
6.	Extremal Properties of Eigenvalues, 104	
7.	Some Additional Results Concerning Eigenvalues, 111	
	Problems, 122	
<b>4.</b>	<b>Matrix Factorizations and Matrix Norms</b>	<b>131</b>
1.	Introduction, 131	
2.	The Singular Value Decomposition, 131	
3.	The Spectral Decomposition and Square Root Matrices of a Symmetric Matrix, 138	
4.	The Diagonalization of a Square Matrix, 144	
5.	The Jordan Decomposition, 147	
6.	The Schur Decomposition, 149	
7.	The Simultaneous Diagonalization of Two Symmetric Matrices, 154	
8.	Matrix Norms, 157	
	Problems, 162	
<b>5.</b>	<b>Generalized Inverses</b>	<b>170</b>
1.	Introduction, 170	
2.	The Moore–Penrose Generalized Inverse, 171	
3.	Some Basic Properties of the Moore–Penrose Inverse, 174	
4.	The Moore–Penrose Inverse of a Matrix Product, 180	
5.	The Moore–Penrose Inverse of Partitioned Matrices, 185	
6.	The Moore–Penrose Inverse of a Sum, 186	
7.	The Continuity of the Moore–Penrose Inverse 188	
8.	Some Other Generalized Inverses, 190	
9.	Computing Generalized Inverses, 197	
	Problems, 204	



**6. Systems of Linear Equations**

210

1. Introduction, 210
2. Consistency of a System of Equations, 210
3. Solutions to a Consistent System of Equations, 213
4. Homogeneous Systems of Equations, 219
5. Least Squares Solutions to a System of Linear Equations, 222
6. Least Squares Estimation For Less Than Full Rank Models, 228
7. Systems of Linear Equations and the Singular Value Decomposition, 233
8. Sparse Linear Systems of Equations, 235  
Problems, 241

**7. Special Matrices and Matrix Operators**

247

1. Introduction, 247
2. Partitioned Matrices, 247
3. The Kronecker Product, 253
4. The Direct Sum, 260
5. The Vec Operator, 261
6. The Hadamard Product, 266
7. The Commutation Matrix, 276
8. Some Other Matrices Associated with the Vec Operator, 283
9. Nonnegative Matrices, 288
10. Circulant and Toeplitz Matrices, 300
11. Hadamard and Vandermonde Matrices, 305  
Problems, 309

**8. Matrix Derivatives and Related Topics**

323

1. Introduction, 323
2. Multivariable Differential Calculus, 323
3. Vector and Matrix Functions, 326
4. Some Useful Matrix Derivatives, 332
5. Derivatives of Functions of Patterned Matrices, 335
6. The Perturbation Method, 337
7. Maxima and Minima, 344
8. Convex and Concave Functions, 349
9. The Method of Lagrange Multipliers, 353  
Problems, 360

<b>9. Some Special Topics Related to Quadratic Forms</b>	<b>370</b>
1. Introduction, 370	
2. Some Results on Idempotent Matrices, 370	
3. Cochran's Theorem, 374	
4. Distribution of Quadratic Forms in Normal Variates, 378	
5. Independence of Quadratic Forms, 384	
6. Expected Values of Quadratic Forms, 390	
7. The Wishart Distribution, 398	
Problems, 409	
<b>References</b>	<b>416</b>
<b>Index</b>	<b>421</b>
<b>List of Series Titles</b>	

# 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 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 that covers some basic statistical concepts and by having most of the statistical examples deal with applications that 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 master's 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 that 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*

## CHAPTER ONE

# A Review of Elementary Matrix Algebra

### 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.

### 2. DEFINITIONS AND NOTATION

Except when stated otherwise, a scalar such as  $\alpha$  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 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  $m$ . An  $m \times 1$  matrix

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

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  $A$  sometimes will be referred to by  $(A)_i$  and  $(A)_j$ , respectively. We will usually use capital letters to represent matrices and lowercase bold letters for vectors.

The diagonal element 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  $A = \text{diag}(a_1, \dots, a_m)$  and  $b$  is a scalar, then we will use  $A^b$  to denote the diagonal matrix  $\text{diag}(a_1^b, \dots, a_m^b)$ . For any  $m \times m$  matrix  $A$ ,  $D_A$  will denote the diagonal matrix with diagonal elements equal to the diagonal elements 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 which has all of its elements below the diagonal equal to 0, while a lower triangular matrix has all of its elements above the diagonal equal to 0.

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 which 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, while 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  $\mathbf{1}_m$  or simply  $\mathbf{1}$  when the size of the vector is obvious.

### 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  $\alpha$  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

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

There is a similar definition for  $BA$ , the postmultiplication of  $B$  by  $A$ . 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 are easy to verify.

**Theorem 1.1.** Let  $\alpha$  and  $\beta$  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)  $(AB)C = A(BC)$ .

#### 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

$$((AB)')_{ij} = (AB)_{ji} = (A)_{j \cdot} (B)_{\cdot i} = \sum_{k=1}^p a_{jk} b_{ki} = (B')_{i \cdot} (A')_{\cdot j} = (B'A')_{ij}$$

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

**Theorem 1.2.** Let  $\alpha$  and  $\beta$  be scalars and  $A$  and  $B$  be matrices. Then, when defined, the following 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, while  $A$  is called a skew-symmetric matrix 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 the previous section 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

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

## 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 the following theorem.

**Theorem 1.3.** Let  $\alpha$  be a scalar and  $A$  and  $B$  be matrices. Then, when the appropriate operations are defined, we have

(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  $A = (0)$ .

## 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  $(1, \dots, m)$ . 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  $A$ .

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  $A_{ij} = (-1)^{i+j} m_{ij}$ . For any  $i = 1, \dots, m$ , the determinant of  $A$  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  $j$ th column,

$$|A| = \sum_{i=1}^m a_{ij} A_{ij} \quad (1.2)$$

On the other hand, if the cofactors of a row or column are matched with the elements from a different row or column, the expansion reduces to 0; that is, if  $k \neq i$ , 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  $5 \times 5$  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  $4 \times 4$  matrix, we get

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

Since the determinant of the  $3 \times 3$  matrix above is 6 we have

$$|A| = 2 \cdot (-1) \cdot 6 = -12$$

The following properties of the determinant are fairly straightforward to verify using the definition of a determinant or the expansion formulas given in (1.1) and (1.2).

**Theorem 1.4.** If  $\alpha$  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) If two rows (or columns) of  $A$  are proportional to one another,  $|A| = 0$ .
- (f) The interchange of two rows (or columns) of  $A$  changes the sign of  $|A|$ .
- (g) If all the elements of a row (or column) of  $A$  are multiplied by  $\alpha$ , then the determinant is multiplied by  $\alpha$ .

- (h) The determinant of  $A$  is unchanged when a multiple of one row (or column) is added to another row (or column).

Consider the  $m \times m$  matrix  $C$  whose columns are given by the vectors  $c_1, \dots, c_m$ ; that is, we can write  $C = (c_1, \dots, c_m)$ . Suppose that, for some vector  $b = (b_1, \dots, b_m)'$  and matrix  $A = (a_1, \dots, a_m)$ , we have

$$c_i = Ab = \sum_{i=1}^m b_i a_i$$

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

$$\begin{aligned} |C| &= \sum_{j=1}^m c_{j1} C_{j1} = \sum_{j=1}^m \left( \sum_{i=1}^m b_i a_{ji} \right) C_{j1} \\ &= \sum_{i=1}^m b_i \left( \sum_{j=1}^m a_{ji} C_{j1} \right) = \sum_{i=1}^m b_i |(a_i, c_2, \dots, c_m)|, \end{aligned}$$

*matrix X* ↙

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 derivation above on each column of  $C$  we find that

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

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

$$|(a_{i_1}, \dots, a_{i_m})| = 0$$

if  $i_j = i_k$  for any  $j \neq k$ . Finally, reordering the columns in  $|(a_{i_1}, \dots, a_{i_m})|$  and

using Theorem 1.4(f), we have

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

This very useful result is summarized below.

**Theorem 1.5.** If both  $A$  and  $B$  are square matrices of the same order, then

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

## 7. THE INVERSE

An  $m \times m$  matrix  $A$  for which  $|A| \neq 0$  is said to be a nonsingular matrix. In this case, there exists a nonsingular matrix denoted by  $A^{-1}$  and called the inverse of  $A$ , such that

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

This inverse is unique since, 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 can be easily verified by utilizing (1.4).

**Theorem 1.6.** If  $\alpha$  is a nonzero scalar, and  $A$  and  $B$  are nonsingular  $m \times m$  matrices, then

- (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

## THE INVERSE

of  $a_{ji}$ . Then

$$\underline{AA\# = A\#A = \text{diag}(|A|, \dots, |A|) = |A|I_m,}$$

since  $(A)_{i \cdot} (A\#)_{\cdot i} = (A\#)_{i \cdot} (A)_{\cdot i} = |A|$  follows 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

$$\begin{aligned} A^{-1}A &= |A|^{-1}|A|I_m \\ &= |A|^{-1}A\#A \\ \underline{A^{-1} &= |A|^{-1}A\#} \quad \leftarrow \end{aligned}$$

The relationship between the inverse of a matrix product and the product of the inverses, given in Theorem 1.6(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 the following result, which is sometimes useful.

**Theorem 1.7.** 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

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

If  $m = n$  and  $C$  and  $D$  are identity matrices, then we obtain the following special case of Theorem 1.7.

**Corollary 1.7.1.** Suppose that  $A$ ,  $B$  and  $A + B$  are all  $m \times m$  nonsingular matrices. Then

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

We obtain another special case of Theorem 1.7 when  $n = 1$ .

**Corollary 1.7.2.** Let  $A$  be an  $m \times m$  nonsingular matrix. If  $c$  and  $d$  are both  $m \times 1$  vectors and  $A + cd'$  is nonsingular, then

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

**Example 1.2.** Theorem 1.7 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  $5 \times 5$  matrix directly. However, the calculations in Theorem 1.7 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

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

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

## 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  $p_1 + p_2 = p$ . 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 the  $2 \times 2$  partitioning given above. For instance, we could partition only the columns of  $A$ , yielding the expression

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

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

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

The matrix  $A$  above 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  $5 \times 5$  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

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



As a result, we have

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

## 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. 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  $\text{rank}(A) = 0$ .

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

- The interchange of two rows (or columns) of  $A$ .
- The multiplication of a row (or column) of  $A$  by a nonzero scalar.
- The addition of a scalar multiple of a row (or column) of  $A$  to another row (or column) 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, while 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 the following very useful result.

**Theorem 1.8.** Let  $A$  be an  $m \times n$  matrix,  $B$  an  $m \times m$  matrix, and  $C$  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 to another matrix of simpler form having the same rank as  $A$ .

**Theorem 1.9.** If  $A$  is an  $m \times n$  matrix of rank  $r > 0$ , then there exist nonsingular  $m \times m$  and  $n \times n$  matrices  $B$  and  $C$ , such that  $H = BAC$  and  $A = B^{-1}HC^{-1}$ , where  $H$  is given by

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

The following is an immediate consequence of Theorem 1.9.

**Corollary 1.9.1.** Let  $A$  be an  $m \times n$  matrix with  $\text{rank}(A) = r > 0$ . Then there exist an  $m \times r$  matrix  $F$  and an  $r \times n$  matrix  $G$  such that  $\text{rank}(F) = \text{rank}(G) = r$  and  $A = FG$ .

## 10. ORTHOGONAL MATRICES

An  $m \times 1$  vector  $p$  is said to be a normalized vector or a unit vector if  $p'p = 1$ . The  $m \times 1$  vectors,  $p_1, \dots, p_n$ , where  $n \leq m$ , are said to be orthogonal if  $p'_i p_j = 0$  for all  $i \neq j$ . If in addition, each  $p_i$  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

$$P'P = I$$

Taking the determinant of both sides, we see that

$$|P'P| = |P'| |P| = |P|^2 = |I| = 1$$

Thus,  $|P| = +1$  or  $-1$ , so that  $P$  is nonsingular,  $P^{-1} = P'$ , and  $PP' = I$  in addition to  $P'P = I$ ; 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 the following theorem.

**Theorem 1.10.** 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.

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 the 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  $A$ .

## 11. QUADRATIC FORMS

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

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

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

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

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

$$x' Bx = \frac{1}{2} x'(A + A')x = \frac{1}{2} (x' Ax + x' A' x) = \frac{1}{2} (x' Ax + x' Ax) = x' Ax$$

since  $x' A' x = (x' A' x)' = x' Ax$ .

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

- (a) If  $x' Ax > 0$  for all  $x \neq 0$ , then  $A$  is positive definite.
- (b) If  $x' Ax \geq 0$  for all  $x \neq 0$  and  $x' Ax = 0$  for some  $x \neq 0$ , then  $A$  is positive semidefinite.

- (c) If  $x'Ax < 0$  for all  $x \neq 0$ , then  $A$  is negative definite.
- (d) If  $x'Ax \leq 0$  for all  $x \neq 0$  and  $x'Ax = 0$  for some  $x \neq 0$ , then  $A$  is negative semidefinite.
- (e) If  $x'Ax > 0$  for some  $x$  and  $x'Ax < 0$  for some  $x$ , 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 9, we will develop some of the most important results involving quadratic forms that are of particular interest in statistics.

## 12. COMPLEX MATRICES

Throughout this entire 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  $c$ , while  $b$  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

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

while 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  $(r, \theta)$ , where  $r$  is the length of the line from the origin to the point  $(a, b)$  and  $\theta$  is the angle between this line and the positive half of the real axis. The relationship between  $a$  and  $b$  and  $r$  and  $\theta$  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 since  $a^2 + b^2 = r^2$  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 the identity above 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 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 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  $C^*C = I_m$ . This is the generalization of the concept of orthogonal matrices to complex matrices since if  $C$  is real then  $C^* = C'$ .

### 13. RANDOM VECTORS AND SOME RELATED STATISTICAL CONCEPTS

In this section, we review some of the basic definitions and results in distribution theory which 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 (1990) or Lindgren (1993). To be consistent with our notation, which uses 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  $\mathcal{B}$  is a subset of the real line, then

$$P(x \in \mathcal{B}) = \int_{\mathcal{B}} f_x(t) dt$$

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  $x$  and  $y$  are random variables and  $\alpha$  and  $\beta$  are constants, then the expectation operator satisfies the properties

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

and

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

where  $g_1$  and  $g_2$  are any real-valued functions. The set of expected values of a random variable  $x$  given by  $E(x^k)$ ,  $k = 1, 2, \dots$ , are known as the moments of  $x$ . These 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  $\sigma_x^2$  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

$\mu_x$ . 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 imbedded 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 since

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

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

We now focus on two particular families of distributions that we will encounter later in this text. A random variable  $x$  is said to have a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , 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  $z$  which 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 a later chapter, are 0, 1, 0, 3, 0, and 15, respectively.

If  $r$  is a positive integer, then a random variable  $v$  has a chi-squared distri-



bution 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 since 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  $\lambda$ . 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  $\chi_r^2$  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,  $p_x(t)$  if all of the random variables are discrete, and a multivariate density function  $f_x(t)$ , if all of the random variables are continuous, where  $x = (x_1, \dots, x_m)'$  and  $t = (t_1, \dots, t_m)'$ . For instance, if they are continuous and  $\mathcal{B}$  is a region in  $R^m$ , then the probability that  $x$  falls in  $\mathcal{B}$  is

$$P(x \in \mathcal{B}) = \int_{\mathcal{B}} \cdots \int f_x(t) dt_1 \cdots dt_m,$$

while the expected value of the real-valued function  $g(x)$  of  $x$  is given by

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

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

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

for all  $t$ .

The mean vector of  $x$ , denoted by  $\mu$ , is the vector of expected values of the  $x_i$ s; that is,

$$\mu = (\mu_1, \dots, \mu_m)' = E(x) = [E(x_1), \dots, E(x_m)]'$$

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  $\sigma_{ij}$  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  $\text{cov}(x_i, x_j) = 0$  since  $E(x_i x_j) = \mu_i \mu_j$ . If  $\alpha_1, \alpha_2, \beta_1$ , and  $\beta_2$  are scalars, 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  $\Omega$ , which has  $\sigma_{ij}$  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  $\Omega$  is a symmetric matrix. Using (1.7) we obtain the matrix formulation for  $\Omega$ ,

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

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

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

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

$$\underline{\text{cov}(y, w) = \text{cov}(\boldsymbol{\alpha}'\mathbf{x}, \boldsymbol{\beta}'\mathbf{x}) = \text{cov}\left(\sum_{i=1}^m \alpha_i x_i, \sum_{j=1}^m \beta_j x_j\right)}$$

$$\underline{= \sum_{i=1}^m \sum_{j=1}^m \alpha_i \beta_j \text{cov}(x_i, x_j) = \sum_{i=1}^m \sum_{j=1}^m \alpha_i \beta_j \sigma_{ij} = \boldsymbol{\alpha}'\Omega\boldsymbol{\beta}}$$

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

$$\underline{E(y) = E(A\mathbf{x}) = AE(\mathbf{x}) = A\boldsymbol{\mu},} \quad (1.8)$$

$$\underline{\text{var}(y) = E\{[\mathbf{y} - E(\mathbf{y})][\mathbf{y} - E(\mathbf{y})]'\} = E[(A\mathbf{x} - A\boldsymbol{\mu})(A\mathbf{x} - A\boldsymbol{\mu})']}$$

$$\underline{= E[A(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})'A'] = A\{E[(\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})']\}A' = A\Omega A'} \quad (1.9)$$

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

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

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

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

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 coefficient  $\rho_{ij}$ , defined by

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

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

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

For any  $m \times 1$  vector  $\alpha$ , 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  $\Omega$  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  $\mu$  and variance  $\sigma^2$ . These quantities can be estimated by the sample mean and sample variance given by

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

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

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})'$$

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  $D_S^{-1/2} = \text{diag}(s_{11}^{-1/2}, \dots, s_{mm}^{-1/2})$ , the correlation matrix can be estimated by the sample correlation matrix defined as

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

The 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  $z = (z_1, \dots, z_m)'$ . The density function of  $z$  is then given by

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

Since  $E(z) = \mathbf{0}$  and  $\text{var}(z) = I_m$ , this particular  $m$ -dimensional multivariate normal distribution is denoted as  $N_m(\mathbf{0}, I_m)$ . If  $\mu$  is an  $m \times 1$  vector of constants and  $T$  is an  $m \times m$  nonsingular matrix, then  $x = \mu + Tz$  is said to have the  $m$ -dimensional multivariate normal distribution with mean vector  $\mu$  and covariance matrix  $\Omega = TT'$ . This is indicated by  $x \sim N_m(\mu, \Omega)$ . For instance, if  $m = 2$ , the vector  $x = (x_1, x_2)'$  has a bivariate normal distribution and its density, induced by the transformation  $x = \mu + Tz$ , 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 rather cumbersome looking density function given in (1.11) can be more conveniently expressed by utilizing 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  $\Omega$  is positive semidefinite, then  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$  is said to have a singular normal distribution. In this case  $\Omega^{-1}$  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  $\boldsymbol{\mu} + Uz$ , 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  $\mathbf{y} \sim N_p(A\boldsymbol{\mu}, A\Omega A')$ .

One of the most widely used procedures in statistics is regression analysis. We will briefly describe this analysis here and subsequently use regression analysis to illustrate some of the matrix methods developed in this text. Some good references on regression are Neter, Wasserman, and Kutner (1985) 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, while 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

$$y = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k + \epsilon, \quad (1.14)$$

where  $\beta_0, \dots, \beta_k$  are unknown parameters and  $\epsilon$  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 being independent random variables with common variance  $\sigma^2$ ; that is, if  $\epsilon_i$  and  $\epsilon_j$  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 (1.14) is an example of a linear model since 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  $x_2 = x_1^2$ . Since 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 (1.14) holds, then we should have, approximately,

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

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

$$y = X\beta$$

if we define

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

One method of estimating the  $\beta_j$ s, which we will discuss from time to time in this text, is called the method of least squares. If  $\hat{\beta} = (\hat{\beta}_1, \dots, \hat{\beta}_k)'$  is an estimate of the parameter vector  $\beta$ , then  $\hat{y} = X\hat{\beta}$  is the vector of fitted values, while  $\hat{y} - y$  gives the vector of errors or deviations of the actual responses from the corresponding fitted values, and

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

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,

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

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  $\beta$  is unique and is given by

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

## PROBLEMS

1. 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  $A = (0)$ .
2. Show that if  $x$  and  $y$  are  $m \times 1$  vectors,  $\text{tr}(xy') = x'y$ . Show that if  $A$  and  $B$  are  $m \times m$  matrices and  $B$  is nonsingular,  $\text{tr}(BAB^{-1}) = \text{tr}(A)$ .
3. Prove Theorem 1.4.
4. Show that any square matrix can be written as the sum of a symmetric matrix and a skew-symmetric matrix.
5. Define the  $m \times m$  matrices,  $A$ ,  $B$ , and  $C$  as

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

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

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

6. Verify the results of Theorem 1.6.
7. Consider the  $4 \times 4$  matrix



## PROBLEMS

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

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

8. Using the matrix  $A$  from the previous problem, verify equation (1.3) when  $i = 1$  and  $k = 2$ .
9. Let  $\lambda$  be a variable and consider the determinant of  $A - \lambda I_m$ , where  $A$  is an  $m \times m$  matrix, as a function of  $\lambda$ . What type of function of  $\lambda$  is this?
10. Find the adjoint matrix of the matrix  $A$  given in Problem 7. Use this to obtain the inverse of  $A$ .
11. Using elementary transformations, determine matrices  $B$  and  $C$  so that  $BAC = I_4$  for the matrix  $A$  given in Problem 7. 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  $A^{-1}$ .
12. Show that the determinant of a triangular matrix is the product of its diagonal elements. In addition, show that the inverse of a lower triangular matrix is a lower triangular matrix.
13. Let  $\mathbf{a}$  and  $\mathbf{b}$  be  $m \times 1$  vectors and  $D$  be an  $m \times m$  diagonal matrix. Use Corollary 1.7.2 to find an expression for the inverse of  $D + \alpha \mathbf{a} \mathbf{b}'$ , where  $\alpha$  is a scalar.
14. 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  $A_{21}$ .

15. 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.

16. Find the rank of the  $4 \times 4$  matrix

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

17. Use elementary transformations to transform the matrix  $A$  given in problem 16 to a matrix  $H$  having the form given in Theorem 1.9. Consequently, determine matrices  $B$  and  $C$  so  $BAC = H$ .
18. Prove parts (b) and (c) of Theorem 1.10.
19. List all permutation matrices of order 3.
20. Consider the  $3 \times 3$  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?

21. 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  $P_1P_1' + P_2P_2' = I_m$ .  $P'P = \begin{bmatrix} P_1' \\ P_2' \end{bmatrix} [P_1 \ P_2] = \begin{bmatrix} I_{m_1} & \\ & I_{m_2} \end{bmatrix}$
22. Let  $A$  be an  $m \times m$  matrix and suppose there exists a real  $n \times m$  matrix  $T$  such that  $T'T = A$ . Show that  $A$  must be nonnegative definite.  
 $\lambda'A\lambda = (T\lambda)'(T\lambda) \geq 0$  □
23. 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.
24. 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.

25. 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.
26. Use properties of expectation to show that for random variables  $x_1$  and  $x_2$ , and scalars  $\alpha_1$ ,  $\alpha_2$ ,  $\beta_1$ , and  $\beta_2$

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

27. Suppose  $\mathbf{x} \sim N_3(\boldsymbol{\mu}, \boldsymbol{\Omega})$ , where

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

and let the  $3 \times 3$  matrix  $A$  and  $2 \times 3$  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  $u = \mathbf{1}'\mathbf{x}$ .
- (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?

28. Suppose  $\mathbf{x}$  is an  $m \times 1$  random vector with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Omega}$ . 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)  $\text{var}[A(\mathbf{x} + \mathbf{c})]$ .

## CHAPTER TWO

# Vector Spaces

### 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. 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  $x_1 \in S$  and  $x_2 \in S$ , then  $x_1 + x_2 \in S$ .
- (b) If  $x \in S$  and  $\alpha$  is any real scalar, then  $\alpha x \in S$ .

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  $T$ . This 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  $x_1 \in S$ ,

## DEFINITIONS

$x_2 \in S$ , and  $\alpha_1$  and  $\alpha_2$  are any real scalars, then  $(\alpha_1 x_1 + \alpha_2 x_2) \in S$ . This can be easily generalized to more than two, say  $n$ , vectors; that is, if  $\alpha_1, \dots, \alpha_n$  are real scalars and  $x_1, \dots, x_n$  vectors such that  $x_i \in S$ , for all  $i$ , then for  $S$  to be a vector space we must have

$$\sum_{i=1}^n \alpha_i x_i \in S \quad (2.1)$$

The left-hand side of (2.1) is called a linear combination of the vectors  $x_1, \dots, x_n$ . Since a vector space is closed under the formation of linear combinations, vector spaces are sometimes also referred to as linear spaces.

**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  $x_1 = (a_1, 0, a_1)'$  and  $x_2 = (a_2, 0, a_2)'$ , where  $a_1$  and  $a_2$  are arbitrary scalars. Then  $x_1 \in S_1$ ,  $x_2 \in S_1$ , and

$$\alpha_1 x_1 + \alpha_2 x_2 = (\alpha_1 a_1 + \alpha_2 a_2, 0, \alpha_1 a_1 + \alpha_2 a_2)' \in S_1,$$

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 the 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 since, for example, if we take  $\alpha = -1$  and  $x = (1, 1, 1)'$ , then  $x \in S_3$  but

$$\alpha x = -(1, 1, 1)' \notin S_3$$

Every vector space with the exception of the vector space  $\{0\}$  has infinitely many vectors. However, by utilizing 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  $S$ .

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  $S$  since 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  $W$ . Since  $u$  and  $v$  are in  $W$ , by the definition of  $W$ , there must exist scalars  $c_1, \dots, c_n$  and  $d_1, \dots, d_n$  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 utilize a product of vectors called an inner product.

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

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

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

$$\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle$$

*Proof.* The result is trivial if  $\mathbf{x} = \mathbf{0}$  since it is easily shown that, in this case,  $\langle \mathbf{x}, \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle = 0$ . 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  $t$

$$0 \leq \langle t\mathbf{x} + \mathbf{y}, t\mathbf{x} + \mathbf{y} \rangle = \langle \mathbf{x}, \mathbf{x} \rangle t^2 + 2\langle \mathbf{x}, \mathbf{y} \rangle t + \langle \mathbf{y}, \mathbf{y} \rangle = at^2 + bt + c$$

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  $\langle \mathbf{x}, \mathbf{y} \rangle^2 \leq \langle \mathbf{x}, \mathbf{x} \rangle \langle \mathbf{y}, \mathbf{y} \rangle$ . □

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  $\mathbf{y}$ .

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  $d(x, y)$  is a distance function defined on the vector space  $S$  if for any vectors  $x, y$ , and  $z$  in  $S$ , we have

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

Property (d) given in the two definitions above 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,  $\langle x, y \rangle$ , the functions,  $\|x\| = \langle x, x \rangle^{1/2}$  and  $d(x, y) = \langle x - y, x - y \rangle^{1/2}$  satisfy the conditions given in Definitions 2.4 and 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  $d_1(x, y) = \|x - y\|_2$ ;  $\|x\|_2$  is the Euclidean norm given by

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

and based on the Euclidean inner product  $\langle x, y \rangle = x'y$ . 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. There are, however, many situations in statistics in which non-Euclidean distance functions are appropriate.

**Example 2.2.** Suppose we wish to compute the distance between the  $m \times 1$  vectors  $x$  and  $\mu$ , where  $x$  is an observation from a distribution having mean vector  $\mu$  and covariance matrix  $\Omega$ . If we want to take into account the effect of the covariance structure, then the Euclidean distance defined above would not be appropriate unless  $\Omega = I$ . 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  $x$  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

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



and is called the Mahalanobis distance between  $\mathbf{x}$  and  $\boldsymbol{\mu}$ . This is sometimes also referred to as the distance between  $\mathbf{x}$  and  $\boldsymbol{\mu}$  in the metric of  $\Omega$  and is useful in a multivariate statistical procedure known as discriminant analysis [see McLachlan (1992) or Huberty (1994)]. Note that if  $\Omega = \mathbf{I}$  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

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

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 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, the norms defined above 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

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

### 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. In order 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, 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  $\alpha$ , a solution will be given by  $\alpha_1 = \alpha$ ,  $\alpha_2 = 0.5\alpha$ ,  $\alpha_3 = -0.5\alpha$ , and so the vectors are linearly dependent. On the other hand, any pair of the three 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 proofs of the following results are left to the reader.

**Theorem 2.3.** Let  $\{x_1, \dots, x_n\}$  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 dependent, 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, there is only one vector 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.

**Theorem 2.4.** The set of  $m \times 1$  vectors  $\{x_1, \dots, x_n\}$ , 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 since 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  $x_n$ , can be expressed as a linear combination of the others; that is, we can find scalars  $\alpha_1, \dots, \alpha_{n-1}$  such that  $x_n = \alpha_1 x_1 + \dots + \alpha_{n-1} x_{n-1}$ . But this implies that

$$\sum_{i=1}^n \alpha_i x_i = \mathbf{0}, \tag{2.2}$$

if we define  $\alpha_n = -1$ , so the vectors  $x_1, \dots, x_n$  are linearly dependent. Conversely, suppose that the vectors  $x_1, \dots, x_n$  are linearly dependent so that (2.2) holds for some choice of  $\alpha_1, \dots, \alpha_n$ , with at least one of the  $\alpha_i$ s, say  $\alpha_n$ , not equal to zero. Thus, we can solve (2.2) for  $x_n$ , in which case we get

$$x_n = \sum_{i=1}^{n-1} \left( \frac{-\alpha_i}{\alpha_n} \right) x_i,$$

so that  $x_n$  is a linear combination of  $x_1, \dots, x_{n-1}$ . 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 1.9 that there are nonsingular  $m \times m$  matrices  $U$  and  $V = [V_1 \ V_2]$ , 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  $U$  gives 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, since  $U$  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, 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 due to 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.  $\square$

#### 4. 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 2.7.** Let  $\{x_1, \dots, x_n\}$  be a set of  $m \times 1$  vectors in a vector space  $S$ . 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  $\mathbf{0}$ , has a basis. Although a basis for a vector space is not uniquely defined, the number of vectors in a basis is unique, and this is what gives us the dimension of a vector space.

**Definition 2.8.** 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  $\{e_1, \dots, e_m\}$ , where for each  $i$ ,  $e_i$  is defined to be the vector whose only nonzero component is the  $i$ th component, which is one. Now, the linear combination of the  $e_i$ s,

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

will equal  $\mathbf{0}$  only if  $\alpha_1 = \dots = \alpha_m = 0$ , so the vectors  $e_1, \dots, e_m$  are linearly independent. Also, if  $x = (x_1, \dots, x_m)'$  is an arbitrary vector in  $R^m$ , then

$$x = \sum_{i=1}^m x_i e_i,$$

so that  $\{e_1, \dots, e_m\}$  spans  $R^m$ . Thus,  $\{e_1, \dots, e_m\}$  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  $\gamma_i$  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.** 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  $S$ . However, the set  $\{\mathbf{x}_1, \mathbf{x}_2\}$  is linearly independent and  $\mathbf{x}_3 = 2\mathbf{x}_1 + \mathbf{x}_2$ , so that  $\{\mathbf{x}_1, \mathbf{x}_2\}$  and  $\{\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3\}$  must span the same vector space. Thus,  $\{\mathbf{x}_1, \mathbf{x}_2\}$  is a basis for  $S$  and so  $S$  is a two-dimensional subspace, that is, a plane in  $R^3$ . Any pair of linearly independent vectors in  $S$  will be a basis for  $S$ ; for example  $\{\mathbf{x}_1, \mathbf{x}_3\}$  and  $\{\mathbf{x}_2, \mathbf{x}_3\}$  are also bases for  $S$ .

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, there may be more than one linear combination that yields a particular  $\mathbf{x}$ . Our next result indicates that this is not the case when the spanning set is a basis.

**Theorem 2.7.** 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.* Since  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  spans  $S$  and  $\mathbf{x} \in S$ , there must exist scalars  $\alpha_1, \dots, \alpha_n$  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 there exists another set of scalars  $\beta_1, \dots, \beta_n$  for which

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

But this then implies that

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

Since  $\{\mathbf{x}_1, \dots, \mathbf{x}_n\}$  is a basis, the vectors  $\mathbf{x}_1, \dots, \mathbf{x}_n$  must be linearly independent and so the equation above is satisfied only if  $\alpha_i - \beta_i = 0$  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 below. The proofs are left to the reader.

**Theorem 2.8**

- (a) Any two bases of a vector space  $S$  must have the same number of vectors.
- (b) If  $\{x_1, \dots, x_n\}$  is a set of linearly independent vectors in a vector space  $S$  and the dimension of  $S$  is  $n$ , then  $\{x_1, \dots, x_n\}$  is a basis for  $S$ .
- (c) If the set  $\{x_1, \dots, x_n\}$  spans the vector space  $S$  and the dimension of  $S$  is  $n$ , then the set  $\{x_1, \dots, x_n\}$  must be linearly independent and thus a basis for  $S$ .
- (d) If the vector space  $S$  has dimension  $n$  and the set of linearly independent vectors  $\{x_1, \dots, x_r\}$  is in  $S$ , where  $r < n$ , there are bases for  $S$  which contain this set as a subset.

**5. 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. We define the following terminology appropriate for such situations.

**Definition 2.9.** 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

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

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

A consequence of Theorem 2.5 is that the dimension of the column space of a square matrix, that is, the number of linearly independent column vectors, is identical to the rank of that matrix when it is nonsingular. The following result shows that this connection between the number of linearly independent columns of a matrix and the rank of that matrix always holds.

**Theorem 2.9.** 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  $\text{rank}(X) = r$  since 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 since, 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  $X$ , there exists some  $(m-r) \times r$  matrix  $A$  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  $X_1$  since 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 since, 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 there exists an  $r \times (n-r)$  matrix  $B$  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

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

Since 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  $\text{rank}(X_{11}) = r$ , since 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  $X_1 = [X_{11} \quad X_{11}B]$ .  $\square$

**Example 2.6.** An implication of Theorem 2.9 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  $(a, b, 0)'$ , while 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

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

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

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 our original definition in terms of submatrices. This is evidenced in the proof of the following basic results regarding the rank of a matrix.

**Theorem 2.10.** Let  $A$  be an  $m \times n$  matrix. Then the following 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,  $\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$ .
- (c)  $\text{rank}(A) = \text{rank}(A') = \text{rank}(AA') = \text{rank}(A'A)$ .

*Proof.* Note that we can write

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

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

$$A + B = [A \quad B] \begin{bmatrix} I_n \\ I_n \end{bmatrix}$$

So using property (a) on the right-hand side of the equation above, we find that

$$\text{rank}(A + B) \leq \text{rank}([A \quad B]) \leq \text{rank}(A) + \text{rank}(B),$$

where the final inequality follows from the fact that the number of linearly independent columns of  $[A \quad B]$  cannot exceed the sum of the numbers of linearly independent columns in  $A$  and  $B$ . In proving (c), note that it follows immediately that  $\text{rank}(A) = \text{rank}(A')$ . It will suffice to prove that  $\text{rank}(A) = \text{rank}(A'A)$  since this can then be used on  $A'$  to prove that  $\text{rank}(A') = \text{rank}\{(A')'A'\} = \text{rank}(AA')$ . If  $\text{rank}(A) = r$ , there exists a full column rank  $m \times r$  matrix  $A_1$  such that, after possibly interchanging some of the columns of  $A$ ,  $A = [A_1 \quad A_1C] = A_1[I_r \quad C]$ , where  $C$  is an  $r \times (n - r)$  matrix. As a result, we have

$$A'A = \begin{bmatrix} I_r \\ C' \end{bmatrix} A_1' A_1 [I_r \quad C]$$

Note that

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

and since 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  $\text{rank}(A'A) = \text{rank}(A_1'A_1) = r$ .  $\square$

The next result 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 2.11.** Let  $A$ ,  $B$ , and  $C$  be any matrices for which the partitioned matrices below are defined. Then

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

$$(b) \operatorname{rank} \left( \begin{bmatrix} A & (0) \\ (0) & B \end{bmatrix} \right) = \operatorname{rank} \left( \begin{bmatrix} (0) & B \\ A & (0) \end{bmatrix} \right) = \operatorname{rank}(A) + \operatorname{rank}(B)$$

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

Our next result gives a useful inequality for the rank of the product of three matrices.

**Theorem 2.12.** Let  $A$ ,  $B$ , and  $C$  be  $p \times m$ ,  $m \times n$ , and  $n \times q$  matrices, respectively. Then

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

*Proof.* It follows from Theorem 2.11(c) that

$$\operatorname{rank} \left( \begin{bmatrix} B & BC \\ AB & (0) \end{bmatrix} \right) \geq \operatorname{rank}(AB) + \operatorname{rank}(BC) \quad (2.3)$$

But, since

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

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

$$\operatorname{rank} \left( \begin{bmatrix} B & BC \\ AB & (0) \end{bmatrix} \right) = \operatorname{rank} \left( \begin{bmatrix} B & (0) \\ (0) & -ABC \end{bmatrix} \right) = \operatorname{rank}(B) + \operatorname{rank}(ABC) \quad (2.4)$$

Combining (2.3) and (2.4) we obtain the desired result.  $\square$

A special case of Theorem 2.12 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 2.10(a).

**Corollary 2.12.1.** If  $A$  is an  $m \times n$  matrix and  $B$  an  $n \times p$  matrix, then

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

## 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, 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 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

$$\begin{aligned} y_1 &= x_1, \\ y_2 &= x_2 - \frac{x_2' y_1}{y_1' y_1} y_1, \\ &\vdots \\ y_r &= x_r - \frac{x_r' y_1}{y_1' y_1} y_1 - \dots - \frac{x_r' y_{r-1}}{y_{r-1}' y_{r-1}} y_{r-1}, \end{aligned} \tag{2.5}$$

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 the following result.

**Theorem 2.13.** 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 2.7 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 rather simple form; premultiplication of the equation for  $x$  above 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  $\{x_1, x_2, x_3\}$ , where

$$x_1 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad x_2 = \begin{bmatrix} 1 \\ -2 \\ 1 \\ -2 \end{bmatrix}, \quad x_3 = \begin{bmatrix} 3 \\ 1 \\ 1 \\ -1 \end{bmatrix}$$

The orthogonal  $y_i$ s are given by  $y_1 = (1, 1, 1, 1)'$ ,

$$y_2 = \begin{bmatrix} 1 \\ -2 \\ 1 \\ -2 \end{bmatrix} - \frac{(-2)}{4} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} = \begin{bmatrix} 3/2 \\ -3/2 \\ 3/2 \\ -3/2 \end{bmatrix},$$

and

$$y_3 = \begin{bmatrix} 3 \\ 1 \\ 1 \\ -1 \end{bmatrix} - \frac{(4)}{(4)} \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} - \frac{(6)}{(9)} \begin{bmatrix} 3/2 \\ -3/2 \\ 3/2 \\ -3/2 \end{bmatrix} = \begin{bmatrix} 1 \\ 1 \\ -1 \\ -1 \end{bmatrix}$$

Normalizing these vectors yields the orthonormal basis  $\{z_1, z_2, z_3\}$ , where

$$z_1 = \begin{bmatrix} 1/2 \\ 1/2 \\ 1/2 \\ 1/2 \end{bmatrix}, \quad z_2 = \begin{bmatrix} 1/2 \\ -1/2 \\ 1/2 \\ -1/2 \end{bmatrix}, \quad z_3 = \begin{bmatrix} 1/2 \\ 1/2 \\ -1/2 \\ -1/2 \end{bmatrix}$$

Thus, for any  $x \in S$ ,  $x = \alpha_1 z_1 + \alpha_2 z_2 + \alpha_3 z_3$ , where  $\alpha_i = x'z_i$ . For instance, since  $x'_3 z_1 = 2$ ,  $x'_3 z_2 = 2$ ,  $x'_3 z_3 = 2$ , we have  $x_3 = 2z_1 + 2z_2 + 2z_3$ .

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

**Theorem 2.14.** Let  $\{z_1, \dots, z_r\}$  be an orthonormal basis for some vector subspace,  $S$ , of  $R^m$ . Then each  $x \in R^m$  can be expressed uniquely as

$$x = u + v,$$

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

*Proof.* It follows from Theorem 2.8(d) that we can find vectors  $z_{r+1}, \dots, z_m$  so that the set  $\{z_1, \dots, z_m\}$  is an orthonormal basis for the  $m$ -dimensional Euclidean space  $R^m$ . It also follows from Theorem 2.7 that there is a unique set of scalars  $\alpha_1, \dots, \alpha_m$  such that

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

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  $S$  due to the orthogonality of the vectors  $\mathbf{z}_1, \dots, \mathbf{z}_m$ .  $\square$

The vector  $\mathbf{u}$  in the theorem above 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  $\mathbf{x}$ .

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  $d_1$  is the Euclidean distance function, then  $d_1(\mathbf{x}, \mathbf{u}) \leq d_1(\mathbf{x}, \mathbf{y})$ . This is fairly 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,

$$\begin{aligned} \{d_1(\mathbf{x}, \mathbf{y})\}^2 &= (\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y}) = \{(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})\}'\{(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})\} \\ &= (\mathbf{x} - \mathbf{u})'(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})'(\mathbf{u} - \mathbf{y}) + 2(\mathbf{x} - \mathbf{u})'(\mathbf{u} - \mathbf{y}) \\ &= (\mathbf{x} - \mathbf{u})'(\mathbf{x} - \mathbf{u}) + (\mathbf{u} - \mathbf{y})'(\mathbf{u} - \mathbf{y}) = \{d_1(\mathbf{x}, \mathbf{u})\}^2 + \{d_1(\mathbf{u}, \mathbf{y})\}^2, \end{aligned}$$

from which  $d_1(\mathbf{x}, \mathbf{u}) \leq d_1(\mathbf{x}, \mathbf{y})$  follows since  $\{d_1(\mathbf{u}, \mathbf{y})\}^2 \geq 0$ .

**Example 2.8.** 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 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

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

In Chapter 8 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  $\hat{\beta}$ ,  $\hat{y} = X\hat{\beta}$  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  $x$ . Thus, the point  $\hat{y}$  that minimizes the distance from  $y$  will be given by the orthogonal projection of  $y$  onto this plane spanned by  $\mathbf{1}_N$  and  $x$ . This means that  $y - \hat{y}$  must be orthogonal to both  $\mathbf{1}_N$  and  $x$ . This leads to the two normal equations

$$0 = (y - \hat{y})'\mathbf{1}_N = y'\mathbf{1}_N - \hat{\beta}'X'\mathbf{1}_N = \sum_{i=1}^N y_i - \hat{\beta}_0 N - \hat{\beta}_1 \sum_{i=1}^N x_i,$$

$$0 = (y - \hat{y})'x = y'x - \hat{\beta}'X'x = \sum_{i=1}^N x_i y_i - \hat{\beta}_0 \sum_{i=1}^N x_i - \hat{\beta}_1 \sum_{i=1}^N x_i^2,$$

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  $\beta_0$ . In this case, the vector of fitted values satisfies  $\hat{y} = \hat{\beta}_0 \mathbf{1}_N$ , so for any choice of  $\hat{\beta}_0$ ,  $\hat{y}$  will be given by a point on the line passing through the origin and  $\mathbf{1}_N$ . Thus if  $\hat{y}$  is to minimize the sum of squared errors and hence the distance from  $y$ , then it must be given by the orthogonal projection of  $y$  onto this line. Consequently, we must have

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

or simply

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

The vector  $v$  in Theorem 2.14 is called the component of  $x$  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,  $S^\perp = \{x: x \in R^m \text{ and } x'y = 0 \text{ for all } y \in S\}$ .

**Theorem 2.15.** 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  $x_1 \in S^\perp$  and  $x_2 \in S^\perp$  so that  $x_1'y = x_2'y = 0$  for any  $y \in S$ . Consequently, for any  $y \in S$  and any scalars  $\alpha_1$  and  $\alpha_2$ , we have

$$(\alpha_1 x_1 + \alpha_2 x_2)'y = \alpha_1 x_1'y + \alpha_2 x_2'y = 0,$$

and so  $(\alpha_1 x_1 + \alpha_2 x_2) \in S^\perp$  and thus  $S^\perp$  is a vector space.  $\square$

A consequence of the following theorem 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 2.16.** Suppose  $\{z_1, \dots, z_m\}$  is an orthonormal basis for  $R^m$  and  $\{z_1, \dots, z_r\}$  is an orthonormal basis for the vector subspace  $S$ . Then  $\{z_{r+1}, \dots, z_m\}$  is an orthonormal basis for  $S^\perp$ .

*Proof.* Let  $T$  be the vector space spanned by  $\{z_{r+1}, \dots, z_m\}$ . We must show that this vector space is the same as  $S^\perp$ . If  $x \in T$  and  $y \in S$ , then there exist scalars  $\alpha_1, \dots, \alpha_m$  such that  $y = \alpha_1 z_1 + \dots + \alpha_r z_r$  and  $x = \alpha_{r+1} z_{r+1} + \dots + \alpha_m z_m$ . Due to the orthogonality of the  $z_i$ s,  $x'y = 0$ , so  $x \in S^\perp$  and thus  $T \subseteq S^\perp$ . Conversely, suppose that  $x \in S^\perp$ . Since  $x$  is also in  $R^m$ , there exist scalars  $\alpha_1, \dots, \alpha_m$  such that  $x = \alpha_1 z_1 + \dots + \alpha_m z_m$ . Now if we let  $y = \alpha_1 z_1 + \dots + \alpha_r z_r$ , then  $y \in S$ , and since  $x \in S^\perp$  we must have  $x'y = \alpha_1^2 + \dots + \alpha_r^2 = 0$ . But this can only happen if  $\alpha_1 = \dots = \alpha_r = 0$ , in which case  $x = \alpha_{r+1} z_{r+1} + \dots + \alpha_m z_m$  and so  $x \in T$ . Thus, we also have  $S^\perp \subseteq T$ , and so this establishes that  $T = S^\perp$ .  $\square$

## 7. PROJECTION MATRICES

The orthogonal projection of an  $m \times 1$  vector  $x$  onto a vector space  $S$  can be conveniently expressed in matrix form. Let  $\{z_1, \dots, z_r\}$  be any orthonormal basis for  $S$  while  $\{z_1, \dots, z_m\}$  is an orthonormal basis for  $R^m$ . Suppose  $\alpha_1, \dots, \alpha_m$  are the constants satisfying the relationship

$$x = (\alpha_1 z_1 + \dots + \alpha_r z_r) + (\alpha_{r+1} z_{r+1} + \dots + \alpha_m z_m) = u + v,$$



where  $u$  and  $v$  are as previously defined. Write  $\alpha = (\alpha_1', \alpha_2')$  and  $Z = [Z_1 \ Z_2]$ , where  $\alpha_1 = (\alpha_1, \dots, \alpha_r)'$ ,  $\alpha_2 = (\alpha_{r+1}, \dots, \alpha_m)'$ ,  $Z_1 = (z_1, \dots, z_r)$ , and  $Z_2 = (z_{r+1}, \dots, z_m)$ . Then the expression for  $x$  given above can be written as

$$x = Z\alpha = Z_1\alpha_1 + Z_2\alpha_2;$$

that is,  $u = Z_1\alpha_1$  and  $v = Z_2\alpha_2$ . Due to the orthonormality of the  $z_i$ s, we have  $Z_1'Z_1 = I_r$  and  $Z_1'Z_2 = (0)$ , and so

$$Z_1Z_1'x = Z_1Z_1'Z\alpha = Z_1Z_1'[Z_1 \ Z_2] \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = [Z_1 \ (0)] \begin{bmatrix} \alpha_1 \\ \alpha_2 \end{bmatrix} = Z_1\alpha_1 = u$$

Thus, we have the following result.

**Theorem 2.17.** 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  $x \in R^m$ , the orthogonal projection of  $x$  onto  $S$  is given by  $Z_1Z_1'x$ .

The matrix  $Z_1Z_1'$  appearing in Theorem 2.17 is called the projection matrix for the vector space  $S$  and sometimes will be denoted by  $P_S$ . Similarly,  $Z_2Z_2'$  is the projection matrix for  $S^\perp$  and  $ZZ' = I_m$  is the projection matrix for  $R^m$ . Since  $ZZ' = Z_1Z_1' + Z_2Z_2'$ , we have the simple equation  $Z_2Z_2' = I_m - Z_1Z_1'$  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 2.18.** 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  $Z_1Z_1' = W_1W_1'$ .

*Proof.* Each column of  $W_1$  can be written as a linear combination of the columns of  $Z_1$  since the columns of  $Z_1$  span  $S$  and each column of  $W_1$  is in  $S$ ; that is, there exists an  $r \times r$  matrix  $P$  such that  $W_1 = Z_1P$ . But  $Z_1'Z_1 = W_1'W_1 = I_r$ , since each matrix has orthonormal columns. Thus,

$$I_r = W_1'W_1 = P'Z_1'Z_1P = P'I_rP = P'P,$$

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

$$W_1W_1' = Z_1PP'Z_1' = Z_1I_rZ_1' = Z_1Z_1' \quad \square$$

We will take another look at the Gram–Schmidt orthonormalization procedure, this time utilizing projection matrices. 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 very easy to verify that for  $i = 1, \dots, r-1$ , the vector  $\mathbf{y}_{i+1}$  can be expressed as

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

that is,  $\mathbf{y}_{i+1} = (\mathbf{I}_m - \mathbf{Z}_{(i)} \mathbf{Z}_{(i)}') \mathbf{x}_{i+1}$ , where  $\mathbf{Z}_{(i)} = (\mathbf{z}_1, \dots, \mathbf{z}_i)$ . 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  $\mathbf{X}_1 = (\mathbf{x}_1, \dots, \mathbf{x}_r)$ , the columns of

$$\mathbf{Z}_1 = \mathbf{X}_1 \mathbf{A} \tag{2.6}$$

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

$$\mathbf{Z}_1' \mathbf{Z}_1 = \mathbf{A}' \mathbf{X}_1' \mathbf{X}_1 \mathbf{A} = \mathbf{I}_r$$

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

$$\mathbf{P}_S = \mathbf{Z}_1 \mathbf{Z}_1' = \mathbf{X}_1 \mathbf{A} \mathbf{A}' \mathbf{X}_1' = \mathbf{X}_1 (\mathbf{X}_1' \mathbf{X}_1)^{-1} \mathbf{X}_1' \tag{2.7}$$

Note that the Gram–Schmidt equations given in (2.5) can be written in matrix form as  $\mathbf{Y}_1 = \mathbf{X}_1 \mathbf{T}$ , where  $\mathbf{Y}_1 = (\mathbf{y}_1, \dots, \mathbf{y}_r)$ ,  $\mathbf{X}_1 = (\mathbf{x}_1, \dots, \mathbf{x}_r)$ , and  $\mathbf{T}$  is an  $r \times r$  upper triangular matrix with each diagonal element equal to 1. The normalization to produce  $\mathbf{Z}_1$  can then be written as  $\mathbf{Z}_1 = \mathbf{X}_1 \mathbf{T} \mathbf{D}^{-1}$ , where  $\mathbf{D}$  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  $\mathbf{A} = \mathbf{T} \mathbf{D}^{-1}$  is upper triangular with positive diagonal elements. Thus the Gram–Schmidt orthonormalization is the particular case of equation (2.6) in which the matrix  $\mathbf{A}$  has been chosen to be the upper triangular square root matrix of  $(\mathbf{X}_1' \mathbf{X}_1)^{-1}$  having positive diagonal elements.

**Example 2.9.** 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

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

This, of course, is the same as  $Z_1Z_1'$ , where  $Z_1 = (z_1, z_2, z_3)$  and  $z_1, z_2, z_3$  are the vectors obtained by the Gram-Schmidt orthonormalization in the previous example. 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  $\mathbf{x}$  since  $\mathbf{x} = \mathbf{x}_3 - \mathbf{x}_1 - \mathbf{x}_2 \in S$ . 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  $\{I - X_1(X_1'X_1)^{-1}X_1'\}\mathbf{x} = \mathbf{x} - X_1(X_1'X_1)^{-1}X_1'\mathbf{x} = \mathbf{x} - \mathbf{u} = (\frac{1}{4}, -\frac{1}{4}, -\frac{1}{4}, \frac{1}{4})'$ . This 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 2.14.

**Example 2.10.** We will generalize some of the ideas of Example 2.8 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

$$y = X\beta + \epsilon,$$

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

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

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

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

minimized. Thus  $X\hat{\beta}$  must be the orthogonal projection of  $y$  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

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

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

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

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

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

and so this sum of squares represents the squared distance of the projection of  $y$  onto the orthogonal complement of the column space of  $X$ . Suppose now that  $\beta$  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  $\beta_1$ , and we wish to decide whether or not  $\beta_2 = \mathbf{0}$ . If the columns of  $X_1$  are orthogonal to the columns of  $X_2$ , then  $X_1'X_2 = (0)$  and

$$(X'X)^{-1} = \begin{bmatrix} (X_1'X_1)^{-1} & (0) \\ (0) & (X_2'X_2)^{-1} \end{bmatrix},$$

and so  $\hat{\beta}$  can be partitioned as  $\hat{\beta} = (\hat{\beta}_1', \hat{\beta}_2')'$ , where  $\hat{\beta}_1 = (X_1'X_1)^{-1}X_1'y$  and  $\hat{\beta}_2 = (X_2'X_2)^{-1}X_2'y$ . Further, the sum of squared errors for the fitted model,  $\hat{y} = X\hat{\beta}$ , can be decomposed as

$$\begin{aligned} (y - X\hat{\beta})'(y - X\hat{\beta}) &= y'(I_N - X(X'X)^{-1}X')y \\ &= y'(I_N - X_1(X_1'X_1)^{-1}X_1' - X_2(X_2'X_2)^{-1}X_2')y \end{aligned}$$

On the other hand, the least squares estimator of  $\beta_1$  in the reduced model

$$y = X_1\beta_1 + \epsilon$$

is  $\hat{\beta}_1 = (X_1'X_1)^{-1}X_1'y$ , while its sum of squared errors is given by

$$SSE_2 = (y - X_1\hat{\beta}_1)'(y - X_1\hat{\beta}_1) = y'(I_N - X_1(X_1'X_1)^{-1}X_1')y$$

Thus, the term  $SSE_2 - SSE_1 = y'X_2(X_2'X_2)^{-1}X_2'y$  gives the reduction in the sum of squared errors attributable to the inclusion of the term  $X_2\beta_2$  in the model  $y = X\beta + \epsilon = X_1\beta_1 + X_2\beta_2 + \epsilon$ , and so its relative size will be helpful in deciding whether or not  $\beta_2 = 0$ . If  $\beta_2 = 0$ , 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, while if  $\beta_2 \neq 0$ , 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, since 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, while  $SSE_2 - SSE_1$  gives the sum of squared deviations in  $k_2$  orthogonal directions, where  $k_2$  is the number of components in  $\beta_2$ . Thus,  $SSE_1/(N - k - 1)$  and  $(SSE_2 - SSE_1)/k_2$  should be of similar magnitudes if  $\beta_2 = 0$ , while the latter should be larger than the former if  $\beta_2 \neq 0$ . Consequently, a decision about  $\beta_2$  can be based on the value of the statistic

$$F = \frac{(SSE_2 - SSE_1)/k_2}{SSE_1/(N - k - 1)} \quad (2.8)$$

Using results that we will develop in Chapter 9, it can be shown that  $F \sim F_{k_2, N-k-1}$  if  $\epsilon \sim N_N(0, \sigma^2 I_N)$  and  $\beta_2 = 0$ .

When  $X_1'X_2 \neq 0$ , the expression for  $(SSE_2 - SSE_1)$  is not equal to  $y'X_2(X_2'X_2)^{-1}X_2'y$  since, in this case,  $\hat{y}$  is not the sum of the projection of  $y$  onto the column space of  $X_1$  and the projection of  $y$  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  $\hat{y}$  into the sum of the projection of  $y$  onto the column space of  $X_1$  and the projection of  $y$  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

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

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

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

Further, the sum of squared errors is given by

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

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

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

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

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

Since  $X$  does not have full column rank, the dimension of its null space is at least one, and so we will be able to find a nonnull vector  $a$  satisfying  $Xa = 0$ . In this case,  $\hat{\beta} + a$  is also a least squares estimator since

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

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, since  $Z_1'Z_1 = I_r$ , it is also idempotent; that is, every projection matrix is symmetric and idempotent. Our next

result proves the converse. Every symmetric idempotent matrix is a projection matrix for some vector space.

**Theorem 2.19.** Let  $P$  be an  $m \times m$  symmetric idempotent matrix of rank  $r$ . Then there is an  $r$ -dimensional vector space which has  $P$  as its projection matrix.

*Proof.* From Corollary 1.9.1, there exist an  $m \times r$  matrix  $F$  and an  $r \times m$  matrix  $G$  such that  $\text{rank}(F) = \text{rank}(G) = r$  and  $P = FG$ . Since  $P$  is idempotent, we have

$$FGFG = FG,$$

which implies that

$$F'FGFGG' = F'FGG' \quad (2.9)$$

Since  $F$  and  $G'$  are full column rank, the matrices  $F'F$  and  $GG'$  are nonsingular. Premultiplying (2.9) 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 to equation (2.7), 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 2.11.** Consider the  $3 \times 3$  matrix

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

Clearly,  $P$  is symmetric and it is easily verified that  $P$  is 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 while 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  $S$  since they are linearly independent and the dimension of  $S$  is 2. Consequently,  $S$  contains all vectors of the form  $(5a - b, 5b - a, 2a + 2b)'$ .

## 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  $x \in R^m$ ,  $u = u(x) = P_S x$  is the orthogonal projection of  $x$  onto  $S$ ; that is, each  $x \in R^m$  is transformed into a  $u \in S$ . The function  $u(x) = P_S x$  is an example of a linear transformation of  $R^m$  into  $S$ .

**Definition 2.11.** Let  $u$  be a function defined for all  $x$  in the vector space  $T$  such that for any  $x \in T$ ,  $u = u(x) \in S$ , where  $S$  is also a vector space. Then the transformation defined by  $u$  is a linear transformation of  $T$  into  $S$  if for any two scalars  $\alpha_1$  and  $\alpha_2$  and any two vectors  $x_1 \in T$  and  $x_2 \in T$ ,

$$u(\alpha_1 x_1 + \alpha_2 x_2) = \alpha_1 u(x_1) + \alpha_2 u(x_2)$$

We will be interested in matrix transformations of the form  $u = Ax$ , where  $x$  is in the subspace of  $R^n$  denoted by  $T$ ,  $u$  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 since for scalars  $\alpha_1, \alpha_2$  and  $n \times 1$  vectors  $x_1$  and  $x_2$ , it follows immediately that

$$A(\alpha_1 x_1 + \alpha_2 x_2) = \alpha_1 Ax_1 + \alpha_2 Ax_2 \quad (2.10)$$

In fact, every linear transformation can be expressed as a matrix transformation. 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 2.10, we saw that for any  $N \times 1$  vector of observations  $y$ , the vector of estimated or fitted values was given by  $\hat{y} = X(X'X)^{-1}X'y$ . Thus, since  $y \in R^N$  and  $\hat{y} \in R(X)$ , we have here a linear transformation of  $R^N$  into  $R(X)$ .

It should be obvious from (2.10) that if  $S$  is actually defined to be the set  $\{u: u = Ax; x \in T\}$ , then  $T$  being a vector space guarantees that  $S$  will also be a vector space. In addition, if the vectors  $x_1, \dots, x_r$  span  $T$ , then the vectors  $Ax_1, \dots, Ax_r$  span  $S$ . In particular, if  $T$  is  $R^n$ , then since  $e_1, \dots, e_n$  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  $A$  since it is spanned by the columns of  $A$ .

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

**Theorem 2.20.** Let the linear transformation of  $R^n$  into  $S$  be given by  $u = Ax$ , where  $x \in R^n$  and  $A$  is an  $m \times n$  matrix. Then the null space of  $A$ , given



by the set

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

is a vector space.

*Proof.* Let  $x_1$  and  $x_2$  be in  $N(A)$  so that  $Ax_1 = Ax_2 = \mathbf{0}$ . Then for any scalars  $\alpha_1$  and  $\alpha_2$ , we have

$$A(\alpha_1x_1 + \alpha_2x_2) = \alpha_1Ax_1 + \alpha_2Ax_2 = \alpha_1(\mathbf{0}) + \alpha_2(\mathbf{0}) = \mathbf{0},$$

so that  $(\alpha_1x_1 + \alpha_2x_2) \in N(A)$  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  $A$ . The following result is an immediate consequence of Theorem 2.16.

**Theorem 2.21.** 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

$$\text{rank}(A) = n - \dim\{N(A)\} \quad (2.11)$$

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

**Example 2.12.** To illustrate the utility of (2.11), we will give an alternative proof of the identity  $\text{rank}(A) = \text{rank}(A'A)$ , which was given as Theorem 2.10(c). Suppose  $x$  is in the null space of  $A$  so that  $Ax = \mathbf{0}$ . Then, clearly, we must have  $A'Ax = \mathbf{0}$ , which implies that  $x$  is also in the null space of  $A'A$ , so it follows that  $\dim\{N(A)\} \leq \dim\{N(A'A)\}$ , or equivalently

$$\text{rank}(A) \geq \text{rank}(A'A) \quad (2.12)$$

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

$$\text{rank}(A) \leq \text{rank}(A'A) \quad (2.13)$$

Combining (2.12) and (2.13), 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, since 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 2.7 there exist scalars  $u_1, \dots, u_m$  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 2.13.** 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  $\theta$ , 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  $\theta$ . 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  $\theta$  followed by a reflection of  $\mathbf{e}_3$  about the  $\mathbf{x}_1, \mathbf{x}_2$  plane.

Although orthogonal transformations are very common, there are situations in which nonsingular nonorthogonal transformations are useful.

**Example 2.14.** 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  $\Omega$ . If we are interested in how these vectors differ from one another, then a plot of the points in  $R^3$  may be useful. However, as discussed in Example 2.2, if  $\Omega$  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 a later chapter that since  $\Omega$  is positive definite, there exists a nonsingular matrix  $T$  satisfying  $\Omega = TT'$ . 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

$$\begin{aligned} d_{\Omega}(\mathbf{x}_i, \mathbf{x}_j) &= \{(\mathbf{x}_i - \mathbf{x}_j)' \Omega^{-1} (\mathbf{x}_i - \mathbf{x}_j)\}^{1/2} \\ &= \{(\mathbf{x}_i - \mathbf{x}_j)' T'^{-1} T^{-1} (\mathbf{x}_i - \mathbf{x}_j)\}^{1/2} \\ &= \{(T^{-1}\mathbf{x}_i - T^{-1}\mathbf{x}_j)' (T^{-1}\mathbf{x}_i - T^{-1}\mathbf{x}_j)\}^{1/2} \\ &= \{(\mathbf{u}_i - \mathbf{u}_j)' (\mathbf{u}_i - \mathbf{u}_j)\}^{1/2} = d_1(\mathbf{u}_i, \mathbf{u}_j), \end{aligned}$$

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

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

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 our next two examples, we discuss some transformations that are sometimes useful in regression analysis.

**Example 2.15.** 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 the components of

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

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

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

can be reexpressed as

$$\begin{aligned} \mathbf{y} &= \beta_0\mathbf{1}_N + \{N^{-1}\mathbf{1}_N\mathbf{1}'_N + (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\}\mathbf{X}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon} \\ &= \boldsymbol{\gamma}_0\mathbf{1}_N + \mathbf{V}_1\boldsymbol{\beta}_1 + \boldsymbol{\epsilon} = \mathbf{V}\boldsymbol{\gamma} + \boldsymbol{\epsilon}, \end{aligned}$$

where  $\mathbf{V} = [\mathbf{1}_N, \mathbf{V}_1] = [\mathbf{1}_N, (\mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}'_N)\mathbf{X}_1]$  and  $\boldsymbol{\gamma} = (\boldsymbol{\gamma}_0, \boldsymbol{\beta}'_1)'$   $= (\beta_0 + N^{-1}\mathbf{1}'_N\mathbf{X}_1\boldsymbol{\beta}_1, \boldsymbol{\beta}'_1)'$ . Since the columns of  $\mathbf{V}_1$  are orthogonal to  $\mathbf{1}_N$ , the least squares estimator of  $\boldsymbol{\gamma}$  simplifies to

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

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

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

then  $(N-1)^{-1}\mathbf{V}'_1\mathbf{V}_1 = S_{22}$  and, since  $\mathbf{V}'_1\mathbf{1}_N = \mathbf{0}$ ,

$$(N-1)^{-1}\mathbf{V}'_1\mathbf{y} = (N-1)^{-1}\mathbf{V}'_1(\mathbf{y} - \bar{y}\mathbf{1}_N) = s_{21}$$

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

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

where  $\delta = (\delta_0, \delta_1')$ ,  $Z = (\mathbf{1}_N, Z_1)$ ,  $\delta_0 = \gamma_0$ ,  $Z_1 = V_1 D_{S_{22}}^{-1/2}$  and  $\delta_1 = D_{S_{22}}^{1/2} \beta_1$ . The least squares estimators are  $\hat{\delta}_0 = \bar{y}$  and  $\hat{\delta}_1 = R_{22}^{-1}r_{21}$ , where we have partitioned the correlation matrix  $R$  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

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

can be written as

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

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

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

is the same as that obtained from the first model.

**Example 2.16.** Consider the multiple regression model

$$y = X\beta + \epsilon,$$

where now  $\text{var}(\boldsymbol{\epsilon}) \neq \sigma^2 \mathbf{I}_N$ . In this case, our previous estimator,  $\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$ , is still the least squares estimator of  $\boldsymbol{\beta}$ , but it doesn't possess certain optimality properties, one of which is illustrated later in Example 3.13, that hold when  $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 \mathbf{I}_N$ . In this example, we will consider the situation in which the  $\epsilon_i$ s are still uncorrelated, but their variances are not all the same. Thus,  $\text{var}(\boldsymbol{\epsilon}) = \boldsymbol{\Omega} = \sigma^2 \mathbf{C}$ , where  $\mathbf{C} = \text{diag}(c_1^2, \dots, c_N^2)$  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  $\boldsymbol{\beta}$  is obtained by making a simple transformation so that ordinary least squares regression applies to the transformed model. Define the matrix  $\mathbf{C}^{-1/2} = \text{diag}(c_1^{-1/2}, \dots, c_N^{-1/2})$  and transform the original regression problem by premultiplying the model equation by  $\mathbf{C}^{-1/2}$ ; the new model equation is

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

or, equivalently,

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

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

$$\text{var}(\boldsymbol{\epsilon}_*) = \text{var}(\mathbf{C}^{-1/2}\boldsymbol{\epsilon}) = \mathbf{C}^{-1/2}\text{var}(\boldsymbol{\epsilon})\mathbf{C}^{-1/2} = \mathbf{C}^{-1/2}\{\sigma^2\mathbf{C}\}\mathbf{C}^{-1/2} = \sigma^2\mathbf{I}_N$$

Thus, for the transformed model, ordinary least squares regression applies and so the least squares estimator of  $\boldsymbol{\beta}$  can be expressed as

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

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

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

A common application related to linear transformations is one in which the matrix  $\mathbf{A}$  and vector  $\mathbf{u}$  consist of known constants, while  $\mathbf{x}$  is a vector of variables, and we wish to determine all  $\mathbf{x}$  for which  $\mathbf{A}\mathbf{x} = \mathbf{u}$ ; that is, we want to find the simultaneous solutions  $x_1, \dots, x_n$  to the system of  $m$  equations

$$a_{11}x_1 + \cdots + a_{1n}x_n = u_1$$

$$\vdots$$

$$a_{m1}x_1 + \cdots + a_{mn}x_n = u_m$$

For instance, in Example 2.10, we saw that the least squares estimator of the parameter vector  $\beta$  in the multiple regression model satisfies the equation,  $X\hat{\beta} = X(X'X)^{-1}X'y$ ; that is, here  $A = X$ ,  $u = X(X'X)^{-1}X'y$ , and  $x = \hat{\beta}$ . In general, if  $u = 0$ , then this system of equations is referred to as a homogeneous system, and the set of all solutions to  $Ax = u$ , in this case, is simply given by the null space of  $A$ . Consequently, if  $A$  has full column rank, then  $x = 0$  is the only solution, whereas there are infinitely many solutions if  $A$  has less than full column rank. A nonhomogeneous system of linear equations is one which has  $u \neq 0$ . While a homogeneous system always has at least one solution,  $x = 0$ , a nonhomogeneous system may or may not have any solutions. A system of linear equations that has no solutions is called an inconsistent system of equations, while a system with solutions is referred to as a consistent system. If  $u \neq 0$  and  $Ax = u$  holds for some  $x$ , then  $u$  must be a linear combination of the columns of  $A$ ; that is, the nonhomogeneous system of equations  $Ax = u$  is consistent if and only if  $u$  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, since  $A^{-1}$  exists, we find that the system  $Ax = u$  has a solution that is unique and is given by  $x = A^{-1}u$ . 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. These results regarding the solution of a general system of linear equations will be developed in Chapter 6.

## 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 utilizes 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

$$S_1 \cap S_2 = \{x \in R^m: x \in S_1 \text{ and } x \in S_2\}$$

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. This follows from the fact that if  $x_1$  and  $x_2$  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, since  $S_1$  and  $S_2$  are vector spaces, for any scalars  $\alpha_1$  and  $\alpha_2$ ,  $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2$  will be in  $S_1$  and  $S_2$  and hence also 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 with the smallest possible dimension containing  $S_1 \cup S_2$ .

**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 the following theorem has been left as an exercise.

**Theorem 2.22.** 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 2.17.** 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

$$\begin{aligned}\mathbf{x}_1 &= (1, 0, 0, 1, 0)', \\ \mathbf{x}_2 &= (0, 0, 1, 0, 1)', \\ \mathbf{x}_3 &= (0, 1, 0, 0, 0)', \\ \mathbf{y}_1 &= (1, 0, 0, 1, 1)', \\ \mathbf{y}_2 &= (0, 1, 1, 0, 0)'\end{aligned}$$

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 there are no constants  $\alpha_1, \alpha_2, \alpha_3, \alpha_4$ , except  $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0$ , satisfying  $\alpha_1\mathbf{x}_1 + \alpha_2\mathbf{x}_2 + \alpha_3\mathbf{x}_3 + \alpha_4\mathbf{y}_1 = \mathbf{0}$ . 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 2.22, 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  $x$ s and the  $y$ s will indicate an appropriate vector; that is, a basis for  $S_1 \cap S_2$  is given by the vector  $y_1 + y_2 = (1, 1, 1, 1, 1)'$ , since  $y_1 + y_2 = x_1 + x_2 + x_3$ .

When  $S_1$  and  $S_2$  are such that  $S_1 \cap S_2 = \{0\}$ , 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  $x \in S_1 \oplus S_2$  has a unique representation as  $x = x_1 + x_2$ , where  $x_1 \in S_1$  and  $x_2 \in S_2$ . A further special case is one in which  $S_1$  and  $S_2$  are orthogonal vector spaces; that is, for any  $x_1 \in S_1$  and  $x_2 \in S_2$ , we have  $x_1'x_2 = 0$ . In this case, the unique representation  $x = x_1 + x_2$  for  $x \in S_1 \oplus S_2$  will have the vector  $x_1$  given by the orthogonal projection of  $x$  onto  $S_1$ , while  $x_2$  will be given by the orthogonal projection of  $x$  onto  $S_2$ . For instance, for any vector subspace  $S$  of  $R^m$ ,  $R^m = S \oplus S^\perp$ , and for any  $x \in R^m$ ,

$$x = P_S x + P_{S^\perp} x$$

In general, if a vector space  $S$  is the sum of the  $r$  vector spaces  $S_1, \dots, S_r$ , and  $S_i \cap S_j = \{0\}$  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 2.18.** 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  $i \leq m-1$ , while  $T_m$  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  $R^m = T_1 \oplus \dots \oplus T_m$ , since it is not true that  $T_i \cap T_j = \{0\}$  for all  $i \neq j$ . Thus any  $x = (x_1, \dots, x_m)'$  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

$$x = x_1 e_1 + \dots + x_m e_m,$$

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  $e_1 \in T_1, e_2 \in T_2, \dots, e_m \in T_m$  or by choosing  $e_2 \in T_1, e_3 \in T_2, \dots, e_m \in T_{m-1}, e_1 \in T_m$ . In addition, the sum of the orthogonal projections of  $x$  onto the spaces  $S_1, \dots, S_m$  yields  $x$ , while the sum of the orthogonal projections of  $x$  onto the spaces  $T_1, \dots, T_m$  yields  $2x$ . 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,  $U_i \cap U_j = \{0\}$  if  $i \neq j$ , so  $R^m = U_1 \oplus \dots \oplus U_m$  and each  $x \in R^m$  has a unique decomposition  $x = x_1 + \dots + x_m$  with  $x_i \in U_i$ . However, in this case, since 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  $U_1, \dots, U_m$ .

## 10. 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. A formal definition follows.

**Definition 2.14.** 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 very 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  $\alpha_1$  and  $\alpha_2$ , while for  $S$  to be a convex set, this need only hold when  $\alpha_1$  and  $\alpha_2$  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  $S$ .

The following result indicates that the intersection of convex sets and the sum of convex sets are themselves convex. The proof will be left as an exercise.

**Theorem 2.23.** 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, and
- (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, due to a generalization of Theorem 2.23(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 our next theorem, we see that the convexity of  $S$  guarantees the convexity of  $\bar{S}$ .

**Theorem 2.24.** 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  $B_n = \{x: x \in R^m, x'x \leq n^{-1}\}$  is a convex set, where  $n$  is a positive integer. Consequently, it follows from Theorem 2.23(b) that  $C_n = S + B_n$  is also convex. It also follows from a generalization of the result given in Theorem 2.23(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,  $T = \{x: x \in R^m, a'x = c\}$ , where  $a$  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, there is a hyperplane which partitions  $R^m$  into two parts so that  $S_1$  is contained in one part, while  $S_2$  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  $0$ .

**Theorem 2.25.** Let  $S$  be a nonempty closed convex subset of  $R^m$  and suppose that  $0 \notin S$ . Then there exists an  $m \times 1$  vector  $a$  such that  $a'x > 0$  for all  $x \in S$ .

*Proof.* Let  $a$  be a point in  $S$  satisfying

$$a'a = \inf_{x \in S} x'x,$$

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  $a \in S$  exists. In addition,  $a \neq 0$  since  $0 \notin S$ . Now let  $c$  be an arbitrary scalar,  $x$  any vector in  $S$  except for  $a$ , and consider the vector  $cx + (1 - c)a$ . The squared length of this vector as a function of  $c$  is given by

$$\begin{aligned} f(c) &= \{cx + (1 - c)a\}'\{cx + (1 - c)a\} = \{c(x - a) + a\}'\{c(x - a) + a\} \\ &= c^2(x - a)'(x - a) + 2ca'(x - a) + a'a \end{aligned}$$

Since the second derivative of this quadratic function  $f(c)$  is positive, we find that it has a unique minimum at the point

$$c_* = \frac{a'(x-a)}{(x-a)'(x-a)}$$

Note that since  $S$  is convex,  $x_c = cx + (1-c)a \in S$  when  $0 \leq c \leq 1$ , and so we must have  $x'_c x_c = f(c) \geq f(0) = a'a$  for  $0 \leq c \leq 1$  due to the way  $a$  was defined. But because of the quadratic structure of  $f(c)$ , this implies that  $f(c) > f(0)$  for all  $c > 0$ . In other words,  $c_* \leq 0$ , and this leads to

$$a'(x-a) \geq 0,$$

or

$$a'x \geq a'a > 0 \quad \square$$

A point  $x_*$  is an interior point of  $S$  if there exists a  $\delta > 0$  such that the set  $S_\delta = \{x: x \in R^m, (x-x_*)'(x-x_*) < \delta\}$  is a subset of  $S$ . On the other hand,  $x_*$  is a boundary point of  $S$  if for each  $\delta > 0$ , the set  $S_\delta$  contains at least one point in  $S$  and at least one point not in  $S$ . Our next result shows that the sets  $S$  and  $\bar{S}$  have the same interior points if  $S$  is convex.

**Theorem 2.26.** Suppose that  $S$  is a convex subset of  $R^m$ , while  $T$  is an open subset of  $R^m$ . If  $T \subset \bar{S}$ , then  $T \subset S$ .

*Proof.* Let  $x_*$  be an arbitrary point in  $T$  and define the sets

$$S_* = \{x: x = y - x_*, y \in S\}, \quad T_* = \{x: x = y - x_*, y \in T\}$$

It follows from the conditions of the theorem that  $S_*$  is convex,  $T_*$  is open, and  $T_* \subset \bar{S}_*$ . The proof will be complete if we can show that  $0 \in S_*$  since this will imply that  $x_* \in S$ . Since  $0 \in T_*$  and  $T_*$  is an open set, we can find an  $\epsilon > 0$  such that each of the vectors,  $\epsilon e_1, \dots, \epsilon e_m, -\epsilon \mathbf{1}_m$  are in  $T_*$ . Since these vectors also must be in  $\bar{S}_*$ , we can find sequences,  $x_{i1}, x_{i2}, \dots$ , for  $i = 1, 2, \dots, m+1$ , such that each  $x_{ij} \in S_*$  and  $x_{ij} \rightarrow \epsilon e_i$  for  $i = 1, \dots, m$ , and  $x_{ij} \rightarrow -\epsilon \mathbf{1}_m$  for  $i = m+1$ , as  $j \rightarrow \infty$ . Define the  $m \times m$  matrix  $X_j = (x_{1j}, \dots, x_{mj})$  so that  $X_j \rightarrow \epsilon I_m$ , as  $j \rightarrow \infty$ . It follows that there exists an integer  $N_1$  such that  $X_j$  is nonsingular for all  $j > N_1$ . For  $j > N_1$ , define

$$y_j = X_j^{-1} x_{m+1,j}, \quad (2.14)$$

so that

$$y_j \rightarrow (\epsilon I_m)^{-1}(-\epsilon \mathbf{1}_m) = -\mathbf{1}_m$$

Thus there exists some integer  $N_2 \geq N_1$ , such that for all  $j > N_2$ , all of the components of  $y_j$  are negative. But from (2.14) we have

$$\mathbf{x}_{m+1,j} - X_j y_j = [X_j \quad \mathbf{x}_{m+1,j}] \begin{bmatrix} -y_j \\ 1 \end{bmatrix} = 0$$

This same equation holds if we replace the vector  $(-y_j, 1)'$  by the unit vector  $(y_j y_j + 1)^{-1/2}(-y_j, 1)'$ . 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 since  $S_*$  is convex,  $\mathbf{0} \in S_*$ .  $\square$

The next result is sometimes called the supporting hyperplane theorem. It states that for any boundary point of a convex set  $S$ , there exists a hyperplane passing through that point such that none of the points of  $S$  are on one side of the hyperplane.

**Theorem 2.27.** 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 there exists an  $m \times 1$  vector  $\mathbf{b} \neq \mathbf{0}$  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  $\mathbf{0} \notin \bar{S}_i$  since  $\mathbf{x}_i \notin \bar{S}$ . Thus, since  $\bar{S}_i$  is closed and convex by Theorem 2.24, it follows from Theorem 2.25 that there exists an  $m \times m$  vector  $\mathbf{a}_i$  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 since  $\mathbf{b}'_i \mathbf{b}_i = 1$ , the sequence  $\mathbf{b}_1, \mathbf{b}_2, \dots$ , is a bounded sequence and so it has a convergent subsequence; that is, there are positive integers  $i_1 < i_2 < \dots$ , and some  $m \times 1$  unit vector  $\mathbf{b}$  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  $\mathbf{x} \in S$  since  $\mathbf{b}'_{i_j}(\mathbf{x} - \mathbf{x}_{i_j}) > 0$  for all  $\mathbf{x} \in S$ . This completes the proof.  $\square$

We are now ready to prove the separating hyperplane theorem.

**Theorem 2.28.** Let  $S_1$  and  $S_2$  be convex subsets of  $R^m$  with  $S_1 \cap S_2 = \emptyset$ . Then there exists an  $m \times 1$  vector  $\mathbf{b} \neq \mathbf{0}$  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 since  $S_2$  is convex. Thus from Theorem 2.23 we know that the set

$$S = S_1 + S_2 = \{x: x = x_1 - x_2, x_1 \in S_1, x_2 \in S_2\}$$

is also convex. In addition,  $0 \notin S$  since  $S_1 \cap S_2 = \emptyset$ . Consequently, using Theorem 2.27, we find that there is an  $m \times 1$  vector  $b \neq 0$  for which  $b'x \geq 0$  for all  $x \in S$ . But this implies that  $b'(x_1 - x_2) \geq 0$  for all  $x_1 \in S_1$  and all  $x_2 \in S_2$ , as is required.  $\square$

Suppose that  $f(x)$  is a nonnegative function which 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  $f(x)$  of the  $m \times 1$  vector  $x$  replaces the interval in  $R^1$  by a symmetric convex set in  $R^m$ . This generalization is due to Anderson (1955). For simple applications of the result to probabilities of random vectors, see Problem 2.44.

**Theorem 2.29.** Let  $S$  be a convex subset of  $R^m$ , symmetric about  $0$ , so that if  $x \in S$ ,  $-x \in S$  also. Let  $f(x) \geq 0$  be a function for which  $f(x) = f(-x)$ ,  $S_\alpha = \{x : f(x) \geq \alpha\}$  is convex for any positive  $\alpha$ , and  $\int_S f(x) dx < \infty$ . Then

$$\int_S f(x + cy) dx \geq \int_S f(x + y) dx,$$

for  $0 \leq c \leq 1$  and  $y \in R^m$ .

A more comprehensive discussion of convex sets can be found in Kelly and Weiss (1979), Lay (1982), and Rockafellar (1970), while some applications of the separating hyperplane theorem to statistical decision theory can be found in Ferguson (1967).

## PROBLEMS

- Determine whether each of the following sets of vectors is a vector space.
  - $\{(a, b, a + b, 1)'\} : -\infty < a < \infty, -\infty < b < \infty\}$ .  $\chi$
  - $\{(a, b, c, a + b - 2c)'\} : -\infty < a < \infty, -\infty < b < \infty, -\infty < c < \infty\}$ .  $\checkmark$
  - $\{(a, b, c, 1 - a - b - c)'\} : -\infty < a < \infty, -\infty < b < \infty, -\infty < c < \infty\}$ .  $\chi$

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)  $\{(1, 0, 0, 0)', (0, 1, 1, 0)', (0, 0, 0, 1)'\}$ .

3. Is the vector  $\mathbf{x} = (1, 1, 1, 1)'$  in the vector space  $S$  given in Problem 2? Is the vector  $\mathbf{y} = (4, 1, 1, 3)'$  in  $S$ ?
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  $V$ .
5. Suppose that  $\mathbf{x}$  is a random vector having a distribution with mean vector  $\boldsymbol{\mu}$  and covariance matrix  $\boldsymbol{\Omega}$  given by

$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 1 \end{bmatrix}, \quad \boldsymbol{\Omega} = \begin{bmatrix} 1 & -0.5 \\ -0.5 & 1 \end{bmatrix}$$

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.

6. Show that the functions  $\|\mathbf{x}\|_p$  and  $\|\mathbf{x}\|_\infty$  defined in Section 2.2 are, in fact, vector norms.
7. Prove Theorem 2.3.
8. Show that the set of vectors  $\{(1, 2, 2, 2)', (1, 2, 1, 2)', (1, 1, 1, 1)'\}$  is a linearly independent set.
9. 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.
10. 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)  $\{(2, 0, 1, 1)', (3, 1, 2, 2)', (2, 1, 1, 2)', (2, 1, 2, 1)'\}$ .
11. Prove the results of Theorem 2.8.
12. Prove that if a set of orthogonal vectors does not contain the null vector, it is a linearly independent set.
13. Find a basis for the vector space given in Problem 2. What is the dimension of this vector space? Find a second different basis for this same vector space.
14. Show that the set of vectors  $\{\gamma_1, \dots, \gamma_m\}$ , given in Example 2.4, is a basis for  $R^m$ .
15. 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  $\text{rank}(AB) = \text{rank}(A)$ .
16. Suppose  $A$  and  $B$  are  $m \times n$  matrices. Show that there exists an  $n \times n$  matrix  $C$  satisfying  $AC = B$  if and only if  $R(B) \subseteq R(A)$ .
17. Prove the results of Theorem 2.11.
18. 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 there exist an  $m \times p$  matrix  $F$  and a  $q \times n$  matrix  $G$  such that  $C = FA + BG$ .

19. Let  $A$  be an  $m \times n$  matrix and  $B$  an  $n \times p$  matrix with  $\text{rank}(B) = n$ . Show that  $\text{rank}(A) = \text{rank}(AB)$ .



## PROBLEMS

20. Refer to Examples 2.7 and 2.9. 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  $AA' = (X_1' X_1)^{-1}$ .
21. 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  $S$ .
22. Using equation (2.7), determine the projection matrix for the vector space  $S$  given in Problem 21. Use this to compute the orthogonal projection of  $x = (1, 0, 0, 1)'$  onto  $S$ .
23. 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  $x = (1, 1, 1)'$ .
24. 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  $S$ .
25. Consider the vector space  $S = \{u: u = Ax, x \in R^4\}$ , where  $A$  is the  $4 \times 4$  matrix given by

$$A = \begin{bmatrix} 1 & 2 & 0 & 1 \\ 1 & 1 & 2 & 2 \\ 1 & 0 & 4 & 3 \\ 1 & 3 & -2 & 0 \end{bmatrix}$$

- Determine the dimension of  $S$  and find a basis.
- Determine the dimension of the null space  $N(A)$  and find a basis for it.
- Is the vector  $(3, 5, 2, 4)'$  in  $S$ ?
- Is the vector  $(1, 1, 1, 1)'$  in  $N(A)$ ?

26. Let  $x \in R^n$  and suppose that  $u(x)$  defines a linear transformation of  $R^n$  into  $R^m$ . Using the basis  $\{e_1, \dots, e_n\}$  for  $R^n$  and the  $m \times 1$  vectors  $u(e_1), \dots, u(e_n)$ , prove that there exists an  $m \times n$  matrix  $A$  for which

$$u(x) = Ax,$$

for every  $x \in R^n$ .

27. Let  $T$  be a vector subspace of  $R^n$  and suppose that  $S$  is the subspace of  $R^m$  given by

$$S = \{u(x) : x \in T\},$$

where the transformation defined by  $u$  is linear. Show that there exists an  $m \times n$  matrix  $A$  satisfying

$$u(x) = Ax,$$

for every  $x \in T$ .

28. Let  $T$  be the vector space spanned by the two vectors  $x_1 = (1, 1, 0)'$  and  $x_2 = (0, 1, 1)'$ . Let  $S$  be the vector space defined as  $S = \{u(x) : x \in T\}$ , where the function  $u$  defines a linear transformation satisfying  $u(x_1) = (2, 3, 1)'$  and  $u(x_2) = (2, 5, 3)'$ . Find a matrix  $A$  such that  $u(x) = Ax$ , for all  $x \in T$ .

29. Consider the linear transformation defined by

$$u(x) = \begin{bmatrix} x_1 - \bar{x} \\ x_2 - \bar{x} \\ \vdots \\ x_m - \bar{x} \end{bmatrix},$$

for all  $x \in R^m$ , where  $\bar{x} = (1/m)\sum x_i$ . Find the matrix  $A$  for which  $u(x) = Ax$  and then determine the dimension of the range and null spaces.

30. 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).

31. 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  $\mu_i$  represents the expected value of a response from treatment  $i$ , while the  $\epsilon_{ij}$ s are independent and identically distributed as  $N(0, \sigma^2)$ .

- (a) If we let  $\beta = (\mu_1, \dots, \mu_k)'$ , write the model above in matrix form by defining  $y$ ,  $X$ , and  $\epsilon$  so that  $y = X\beta + \epsilon$ .
- (b) Find the least squares estimator of  $\beta$  and show that the sum of squared

errors for the corresponding fitted model is given by

$$SSE_1 = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2,$$

where

$$\bar{y}_i = \sum_{j=1}^{n_i} y_{ij}/n_i$$

(c) If  $\mu_1 = \dots = \mu_k = \mu$ , then the reduced model

$$y_{ij} = \mu + \epsilon_{ij}$$

holds for all  $i$  and  $j$ . Find the least squares estimator of  $\mu$  and the sum of squared errors  $SSE_2$  for the fitted reduced model. Show that  $SSE_2 - SSE_1$ , referred to as the sum of squares for treatment and denoted SST, can be expressed as

$$SST = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2,$$

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 (2.8) takes the form

$$F = \frac{SST/(k-1)}{SSE_1/(n-k)}$$

32. 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.

## PROBLEMS

- (a) Show that  $S = \{y: y = X\hat{\beta}, A\hat{\beta} = \mathbf{0}\}$  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  $\hat{\beta}$  is given by

$$\hat{\beta} = C(C'X'XC)^{-1}C'X'y$$

33. 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  $R^m$ .
34. 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  $S_1 + S_2 \subseteq T$ .
35. Prove Theorem 2.22.
36. 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  $S_1 + S_2$ .
37. Let  $S_1$  be the vector space spanned by the vectors

$$x_1 = \begin{bmatrix} 3 \\ 1 \\ 3 \\ 1 \end{bmatrix}, \quad x_2 = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad x_3 = \begin{bmatrix} 2 \\ 1 \\ 2 \\ 1 \end{bmatrix},$$

while 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  $S_1 \cap S_2$ .

38. Let  $S_1$  and  $S_2$  be vector subspaces of  $R^m$  with  $\dim(S_1) = r_1$  and  $\dim(S_2) = r_2$ .
- Find expressions in terms of  $m$ ,  $r_1$ , and  $r_2$  for the smallest and largest possible values of  $\dim(S_1 + S_2)$ .
  - Find the smallest and largest possible values of  $\dim(S_1 \cap S_2)$ .
39. 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  $S_1 \cap S_2 = \{0\}$ .
40. Let  $S_1$  be the vector space spanned by  $\{(1, 1, -2, 0)', (2, 0, 1, -3)'\}$ , while  $S_2$  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  $S_1 \oplus S_2$ ? Are  $S_1$  and  $S_2$  orthogonal vector spaces?
41. 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

$$\dim(T) = \dim(S_1) + \dim(S_2)$$

42. The concept of orthogonal projections and their associated projection matrices can be extended to projections that are not orthogonal. In the case of orthogonal projections onto the vector space  $S \subseteq R^m$ , we decompose  $R^m$  as  $R^m = S \oplus S^\perp$ . The projection matrix that projects orthogonally onto  $S$  is the matrix  $P$  satisfying  $Py \in S$  and  $(y - Py) \in S^\perp$  for all  $y \in R^m$  and  $Px = x$  for all  $x \in S$ . If  $S$  is the column space of the full rank matrix  $X$ , then  $S^\perp$  will be the null space of  $X'$ , and the projection matrix described above is given by  $P = X(X'X)^{-1}X'$ . Suppose now that we decompose  $R^m$  as  $R^m = S \oplus T$ , where  $S$  is as before and  $T$  is the null space of the full rank matrix  $Y'$ . Note that  $S$  and  $T$  are not necessarily orthogonal vector spaces. We wish to find a projection matrix  $Q$  satisfying  $Qy \in S$  and  $(y - Qy) \in T$  for all  $y \in R^m$  and  $Qx = x$  for all  $x \in S$ .
- Show that  $Q$  is a projection matrix if and only if it is an idempotent matrix.
  - Show that  $Q$  can be expressed as  $Q = X(Y'X)^{-1}Y'$ .

43. Prove Theorem 2.23.

44. Show that if  $S_1$  and  $S_2$  are convex subsets of  $R^m$ , then  $S_1 \cup S_2$  need not be convex.

45. Show that for any positive scalar  $n$ , the set  $B_n = \{x: x \in R^m, x'x \leq n^{-1}\}$  is convex.

46. For any set  $S \subseteq R^m$ , show that its convex hull  $C(S)$  consists of all convex combinations of the vectors in  $S$ .
47. 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  $S$ .
48. 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  $\alpha$ . Suppose that  $S$  is a convex subset of  $R^m$ , symmetric about  $\mathbf{0}$ .
- Show that  $P(\mathbf{x} + c\mathbf{y} \in S) \geq P(\mathbf{x} + \mathbf{y} \in S)$  for any constant vector  $\mathbf{y} \in S$  and  $0 \leq c \leq 1$ .
  - Show that the inequality in (a) also holds if  $\mathbf{y}$  is an  $m \times 1$  random vector distributed independently of  $\mathbf{x}$ .
  - Show that if  $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$ , its density function satisfies the conditions of this exercise.
  - 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  $P(\mathbf{x} \in S) \leq P(\mathbf{y} \in S)$ .

## CHAPTER THREE

# Eigenvalues and Eigenvectors

### 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. This is a consequence of some of the properties of eigenvalues and eigenvectors that we will develop in this chapter. But before we get to these properties, we must first understand how eigenvalues and eigenvectors are defined and how they are calculated.

### 2. EIGENVALUES, EIGENVECTORS, AND EIGENSPACES

If  $A$  is an  $m \times m$  matrix, then any scalar  $\lambda$  satisfying the equation

$$Ax = \lambda x, \quad (3.1)$$

for some  $m \times 1$  vector  $x \neq 0$ , is called an eigenvalue of  $A$ . The vector  $x$  is called an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda$ , and equation (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

$$(A - \lambda I)x = 0 \quad (3.2)$$

Note that if  $|A - \lambda I| \neq 0$ , then  $(A - \lambda I)^{-1}$  would exist and so premultiplication of equation (3.2) by this inverse would lead to a contradiction of the already stated assumption that  $x \neq 0$ . Thus, any eigenvalue  $\lambda$  must satisfy the determinantal equation

$$|A - \lambda I| = 0,$$



which is known as the characteristic equation of  $A$ . Using the definition of a determinant, we readily observe that the characteristic equation is an  $m$ th degree polynomial in  $\lambda$ ; that is, there are scalars  $\alpha_0, \dots, \alpha_{m-1}$  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, there are  $m$  scalars  $\lambda_1, \dots, \lambda_m$ , which 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  $3 \times 3$  matrix  $A$  given by

$$A = \begin{bmatrix} 5 & -3 & 3 \\ 4 & -2 & 3 \\ 4 & -4 & 5 \end{bmatrix}$$

The characteristic equation of  $A$  is

$$\begin{aligned} |A - \lambda I| &= \begin{vmatrix} 5 - \lambda & -3 & 3 \\ 4 & -2 - \lambda & 3 \\ 4 & -4 & 5 - \lambda \end{vmatrix} \\ &= -(5 - \lambda)^2(2 + \lambda) - 3(4)^2 - 4(3)^2 \\ &\quad + 3(4)(2 + \lambda) + 3(4)(5 - \lambda) + 3(4)(5 - \lambda) \\ &= -\lambda^3 + 8\lambda^2 - 17\lambda + 10 \\ &= -(\lambda - 5)(\lambda - 2)(\lambda - 1) = 0, \end{aligned}$$

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  $Ax = 5x$  for  $x$ , which yields the system of equations

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

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  $x$  having  $x_1 = x_2 = x_3$ , such as the vector  $(1, 1, 1)'$ , is an eigenvector of  $A$  associated with the root 5. In a similar fashion, by solving the equation  $Ax = \lambda x$ , 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  $x$  satisfies (3.1) for a given value of  $\lambda$ , then so will  $(\alpha x)$  for any nonzero scalar  $\alpha$ . Thus, eigenvectors are not uniquely defined unless we impose some scale constraint; for instance, we might only consider eigenvectors,  $x$ , satisfying  $x'x = 1$ . 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, since each of these eigenvectors, when multiplied by  $-1$ , yields another normalized eigenvector.

The following example 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

$$|A - \lambda I| = \begin{vmatrix} 1 - \lambda & 1 \\ -2 & -1 - \lambda \end{vmatrix} = -(1 - \lambda)(1 + \lambda) + 2 = \lambda^2 + 1 = 0,$$

so that the eigenvalues of  $A$  are  $i = \sqrt{-1}$  and  $-i$ . To find an eigenvector corresponding to the root  $i$ , write  $x = (x_1, x_2)' = (y_1 + iz_1, y_2 + iz_2)'$  and solve for  $y_1, y_2, z_1, z_2$  using the equation  $Ax = ix$ . From this we find that for any real scalar  $\alpha \neq 0$ ,  $x = (\alpha + i\alpha, -2\alpha)'$  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  $x = (\alpha - i\alpha, -2\alpha)'$ .

The  $m$  eigenvalues of a matrix  $A$  need not all be different since 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 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  $\lambda$ , along with the trivial vector  $\mathbf{0}$ , is called the eigenspace of  $A$  associated with  $\lambda$ ; 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 root  $\lambda$ , 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  $\alpha$  and  $\beta$

$$A(\alpha\mathbf{x} + \beta\mathbf{y}) = \alpha A\mathbf{x} + \beta A\mathbf{y} = \alpha(\lambda\mathbf{x}) + \beta(\lambda\mathbf{y}) = \lambda(\alpha\mathbf{x} + \beta\mathbf{y})$$

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 this leads to two linearly independent solutions; 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 the preceding example, for each value of  $\lambda$ , we have  $\dim\{S(\lambda)\}$  being equal to the multiplicity of  $\lambda$ . This is not always the case; the following example illustrates that it is possible for  $\dim\{S(\lambda)\}$  to be less than the multiplicity of the eigenvalue  $\lambda$ .

**Example 3.4.** Consider the  $3 \times 3$  matrix given by

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

Since  $|A - \lambda I| = (1 - \lambda)^3$ ,  $A$  has the eigenvalue 1 repeated three times. The eigenvalue–eigenvector equation  $Ax = \lambda x$  yields the three scalar equations

$$\begin{aligned} x_1 + 2x_2 + 3x_3 &= x_1 \\ x_2 &= x_2 \\ 2x_2 + x_3 &= x_3, \end{aligned}$$

which have as a solution only vectors of the form  $x = (a, 0, 0)'$ . 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. SOME BASIC PROPERTIES OF EIGENVALUES AND EIGENVECTORS

In this section, we establish some very useful results regarding eigenvalues. The proofs of the results in our first theorem, 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

- (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) Each of the eigenvalues of  $A$  is either +1 or -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  $\lambda$  to be less than the multiplicity of  $\lambda$ . The following theorem shows that if  $\dim\{S_A(\lambda)\} \neq r$ , where  $r$  denotes the multiplicity of  $\lambda$ , then  $\dim\{S_A(\lambda)\} < r$ .

**Theorem 3.3.** Suppose  $\lambda$  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  $\lambda$  is an eigenvalue of  $A$ , by definition there exists an  $x \neq 0$  satisfying the eigenvalue–eigenvector equation  $Ax = \lambda x$  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  $\lambda$ . Form a nonsingular  $m \times m$  matrix  $X$  which 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  $\lambda$ , 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  $\mu$  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,  $\lambda$  must be an eigenvalue of  $X^{-1}AX$  with multiplicity of at least  $k$ . The result now follows since, from Theorem 3.2(d), the eigenvalues of  $X^{-1}AX$  are the same as those of  $A$ .  $\square$

We now prove the following theorem involving both the eigenvalues and the eigenvectors of a matrix.

**Theorem 3.4.** Let  $\lambda$  be an eigenvalue of the  $m \times m$  matrix  $A$  and let  $x$  be a corresponding eigenvector. Then,

- (a) If  $n$  is an integer  $\geq 1$ ,  $\lambda^n$  is an eigenvalue of  $A^n$  corresponding to the eigenvector  $x$ .
- (b) If  $A$  is nonsingular,  $\lambda^{-1}$  is an eigenvalue of  $A^{-1}$  corresponding to the eigenvector  $x$ .

*Proof.* Part (a) is proven by repeatedly using the relationship  $Ax = \lambda x$ ; that is, we have

$$A^n x = A^{n-1}(Ax) = A^{n-1}(\lambda x) = \lambda A^{n-1} x = \dots = \lambda^n x$$

To prove part (b), premultiply the eigenvalue–eigenvector equation

$$Ax = \lambda x$$

by  $A^{-1}$ , yielding the equation

$$x = \lambda A^{-1} x \tag{3.3}$$

Since  $A$  is nonsingular, we know from Theorem 3.2(b) that  $\lambda \neq 0$ , and so dividing both sides of (3.3) by  $\lambda$  yields

$$A^{-1} x = \lambda^{-1} x,$$

which is the eigenvalue–eigenvector equation for  $A^{-1}$ , with eigenvalue  $\lambda^{-1}$  and eigenvector  $x$ .  $\square$

The determinant and trace of a matrix have very simple and useful relationships with the eigenvalues of that matrix. These relationships are established in the next theorem.

**Theorem 3.5.** Let  $A$  be an  $m \times m$  matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$ . Then

$$(a) \operatorname{tr}(A) = \sum_{i=1}^m \lambda_i,$$

$$(b) |A| = \prod_{i=1}^m \lambda_i.$$

*Proof.* Recall that the characteristic equation,  $|A - \lambda I| = 0$ , can be expressed in the polynomial form

$$(-\lambda)^m + \alpha_{m-1}(-\lambda)^{m-1} + \dots + \alpha_1(-\lambda) + \alpha_0 = 0 \tag{3.4}$$

We will first identify the coefficients  $\alpha_0$  and  $\alpha_{m-1}$ . We can determine  $\alpha_0$  by evaluating the left-hand side of equation (3.4) at  $\lambda = 0$ ; thus,  $\alpha_0 = |A - (0)I| = |A|$ . In order to find  $\alpha_{m-1}$ , 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  $\alpha_{m-1}$  is the coefficient of  $(-\lambda)^{m-1}$ , to evaluate this term we only need to consider the terms in the sum which involve at least  $m - 1$  of the diagonal elements of  $(A - \lambda I)$ . But each term in the sum is the product of  $m$  elements from the matrix  $(A - \lambda I)$ , multiplied by the appropriate sign, with one element

chosen from each row and each column of  $(A - \lambda I)$ . Consequently, the only term in the sum involving at least  $m - 1$  of the diagonal elements of  $(A - \lambda I)$  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  $\alpha_{m-1}$  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 since  $\lambda_1, \dots, \lambda_m$  are the roots to the characteristic equation, which is an  $m$ th degree polynomial in  $\lambda$ , 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

$$|A| = \prod_{i=1}^m \lambda_i, \quad \text{tr}(A) = \sum_{i=1}^m \lambda_i \quad \square$$

The following theorem gives a sufficient condition for a set of eigenvectors to be linearly independent.

**Theorem 3.6.** 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 since  $\mathbf{x}_1$ , 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 since no eigenvector can be the null vector, such that

$$\alpha_1 \mathbf{x}_1 + \cdots + \alpha_{h+1} \mathbf{x}_{h+1} = \mathbf{0}$$

Premultiplying the left-hand side of this equation by  $(A - \lambda_{h+1} I)$ , we find that

$$\begin{aligned}
& \alpha_1(A - \lambda_{h+1}I)x_1 + \cdots + \alpha_{h+1}(A - \lambda_{h+1}I)x_{h+1} \\
&= \alpha_1(Ax_1 - \lambda_{h+1}x_1) + \cdots + \alpha_{h+1}(Ax_{h+1} - \lambda_{h+1}x_{h+1}) \\
&= \alpha_1(\lambda_1 - \lambda_{h+1})x_1 + \cdots + \alpha_h(\lambda_h - \lambda_{h+1})x_h
\end{aligned}$$

also must be equal to 0. But  $x_1, \dots, x_h$  are linearly independent so it follows that

$$\alpha_1(\lambda_1 - \lambda_{h+1}) = \cdots = \alpha_h(\lambda_h - \lambda_{h+1}) = 0$$

We know that at least one of the scalars  $\alpha_1, \dots, \alpha_h$  is not equal to zero, and if, for instance,  $\alpha_i$  is one of these nonzero scalars, we then must have  $\lambda_i = \lambda_{h+1}$ . This contradicts the conditions of the theorem, so the vectors  $x_1, \dots, x_r$  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 3.6 that the matrix  $X = (x_1, \dots, x_m)$ , where  $x_i$  is an eigenvector corresponding to  $\lambda_i$ , is nonsingular. It also follows from the eigenvalue–eigenvector equation  $Ax_i = \lambda_i x_i$  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 to 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.

Clearly, when a matrix is diagonalizable, its rank equals the number of its nonzero eigenvalues, since

$$\text{rank}(A) = \text{rank}(X^{-1}AX) = \text{rank}(\Lambda)$$

follows from Theorem 1.8. 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 3.5.** Consider the  $2 \times 2$  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.



Our final theorem, known as the Cayley–Hamilton Theorem, states that a matrix satisfies its own characteristic equation. A proof of this result can be found in Hammarling (1970).

**Theorem 3.7.** Let  $A$  be an  $m \times m$  matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$ . Then

$$\prod_{i=1}^m (A - \lambda_i I) = (0);$$

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

$$(-A)^m + \alpha_{m-1}(-A)^{m-1} + \dots + \alpha_1(-A) + \alpha_0 I = (0)$$

#### 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 3.8.** Let  $A$  be an  $m \times m$  real symmetric matrix. Then the eigenvalues of  $A$  are real, and corresponding to any eigenvalue there exist eigenvectors that are real.

*Proof.* Let  $\lambda = \alpha + i\beta$  be an eigenvalue of  $A$  and  $\mathbf{x} = \mathbf{y} + iz$  a corresponding eigenvector, where  $i = \sqrt{-1}$ . We will first show that  $\beta = 0$ . Substitution of these expressions for  $\lambda$  and  $\mathbf{x}$  in the eigenvalue–eigenvector equation  $A\mathbf{x} = \lambda\mathbf{x}$  yields

$$A(\mathbf{y} + iz) = (\alpha + i\beta)(\mathbf{y} + iz) \tag{3.5}$$

Premultiplying (3.5) by  $(\mathbf{y} - iz)'$ , we get

$$(\mathbf{y} - iz)'A(\mathbf{y} + iz) = (\alpha + i\beta)(\mathbf{y} - iz)'(\mathbf{y} + iz),$$

which simplifies to

$$\mathbf{y}'A\mathbf{y} + \mathbf{z}'A\mathbf{z} = (\alpha + i\beta)(\mathbf{y}'\mathbf{y} + \mathbf{z}'\mathbf{z}),$$

since  $y'Az = z'Ay$  follows from the symmetry of  $A$ . Now  $x \neq 0$  implies that  $(y'y + z'z) > 0$  and, consequently, we must have  $\beta = 0$  since the left-hand side of the equation above is real. Substituting  $\beta = 0$  in (3.5), we find that

$$Ay + iAz = \alpha y + i\alpha z$$

Thus,  $x = y + iz$  will be an eigenvector of  $A$  corresponding to  $\lambda = \alpha$  as long as  $y$  and  $z$  satisfy  $Ay = \alpha y$ ,  $Az = \alpha z$ , and at least one is not  $0$  so that  $x \neq 0$ . A real eigenvector is then constructed by selecting  $y \neq 0$ , such that  $Ay = \alpha y$  and  $z = 0$ .  $\square$

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 another. We will now show that, if  $A$  is symmetric, we can say a bit more. Suppose that  $x$  and  $y$  are eigenvectors of  $A$  corresponding to the eigenvalues  $\lambda$  and  $\gamma$ , where  $\lambda \neq \gamma$ . Then, since  $A$  is symmetric, it follows that

$$\lambda x'y = (\lambda x)'y = (Ax)'y = x'A'y = x'(Ay) = x'(\gamma y) = \gamma x'y$$

Since  $\lambda \neq \gamma$ , we must have  $x'y = 0$ ; 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 the following result.

**Theorem 3.9.** Let  $A$  be an  $m \times m$  symmetric matrix and let  $x$  be any nonzero  $m \times 1$  vector. Then for some  $r \geq 1$ , the vector space spanned by the vectors  $x, Ax, \dots, A^{r-1}x$ , contains an eigenvector of  $A$ .

*Proof.* Let  $r$  be the smallest integer for which  $x, Ax, \dots, A^r x$  form a linearly dependent set. Then there exist scalars,  $\alpha_0, \dots, \alpha_r$ , not all of which are zero, such that

$$\alpha_0 x + \alpha_1 Ax + \dots + \alpha_r A^r x = (\alpha_0 I_m + \alpha_1 A + \dots + A^r)x = 0,$$

where without loss of generality we have taken  $\alpha_r = 1$ , since the way  $r$  was chosen guarantees that  $\alpha_r$  is not zero. The expression in the parentheses is an  $r$ th-degree matrix polynomial in  $A$ . This can be factored in a fashion similar to the way scalar polynomials are factored; that is, it can be written as

$$(A - \gamma_1 I_m)(A - \gamma_2 I_m) \cdots (A - \gamma_r I_m),$$

where  $\gamma_1, \dots, \gamma_r$  are the roots of the polynomial satisfying  $\alpha_0 = (-1)^r \gamma_1 \cdot \gamma_2 \cdots \gamma_r, \dots, \alpha_{r-1} = -(\gamma_1 + \gamma_2 + \dots + \gamma_r)$ . Let

$$\begin{aligned} y &= (A - \gamma_2 I_m) \cdots (A - \gamma_r I_m)x, \\ &= (-1)^{r-1} \gamma_2 \cdots \gamma_r x + \cdots + A^{r-1}x, \end{aligned}$$

and note that  $y \neq 0$  since, otherwise,  $x, Ax, \dots, A^{r-1}x$  would be a linearly dependent set, contradicting the definition of  $r$ . Thus,  $y$  is in the space spanned by  $x, Ax, \dots, A^{r-1}x$  and

$$(A - \gamma_1 I_m)y = (A - \gamma_1 I_m)(A - \gamma_2 I_m) \cdots (A - \gamma_r I_m)x = 0$$

Consequently,  $y$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\gamma_1$ , and so the proof is complete.  $\square$

**Theorem 3.10.** 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,  $x_1, \dots, x_h$ , where  $1 \leq h < m$ , then we can find another normalized eigenvector  $x_{h+1}$  orthogonal to each of these vectors. Select any vector  $x$  which is orthogonal to each of the vectors  $x_1, \dots, x_h$ . Note that for any positive integer  $k$ ,  $A^k x$  is also orthogonal to  $x_1, \dots, x_h$  since, if  $\lambda_i$  is the eigenvalue corresponding to  $x_i$ , it follows from the symmetry of  $A$  and Theorem 3.4(a) that

$$x_i' A^k x = \{(A^k)' x_i\}' x = (A^k x_i)' x = \lambda_i^k x_i' x = 0$$

From the previous theorem we know that, for some  $r$ , the space spanned by  $x, Ax, \dots, A^{r-1}x$  contains an eigenvector, say  $y$ , of  $A$ . This vector  $y$  also must be orthogonal to  $x_1, \dots, x_h$  since it is from a vector space spanned by a set of vectors orthogonal to  $x_1, \dots, x_h$ . Thus, we can take  $x_{h+1} = (y'y)^{-1/2}y$ . 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  $X = (x_1, \dots, x_m)$ , where  $x_1, \dots, x_m$  are the orthonormal vectors described in the proof, and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ , then the eigenvalue–eigenvector equation  $Ax_i = \lambda_i x_i$  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  $\Lambda$ , which are the diagonal elements of  $\Lambda$ . Thus, if  $\lambda$  is a multiple root of  $A$  with multiplicity  $r > 1$ , then  $r$  of the diagonal elements of  $\Lambda$  are equal to  $\lambda$  and  $r$  of the eigenvectors, say  $x_1, \dots, x_r$ , correspond to this root  $\lambda$ . Consequently, the dimension of the eigenspace of  $A$ ,  $S_A(\lambda)$ , corresponding to  $\lambda$ , is equal to the multiplicity  $r$ . The set of orthonormal eigenvectors corresponding to this root is not unique. Any orthonormal basis for  $S_A(\lambda)$  will be a set of  $r$  orthonormal vectors associated with the eigenvalue  $\lambda$ . For example, if we let  $X_1 = (x_1, \dots, x_r)$  and let  $Q$  be any  $r \times r$  orthogonal matrix, the columns of  $Y_1 = X_1Q$  also form a set of orthonormal eigenvectors corresponding to  $\lambda$ .

**Example 3.6.** 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  $\beta$  in the model  $y = X\beta + \epsilon$

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

will be imprecise since its covariance matrix

$$\begin{aligned} \text{var}(\hat{\beta}) &= (X'X)^{-1}X'\{\text{var}(y)\}X(X'X)^{-1} \\ &= (X'X)^{-1}X'\{\sigma^2I\}X(X'X)^{-1} = \sigma^2(X'X)^{-1} \end{aligned}$$

will tend to have some large elements due to 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  $x_l$ , 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

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

discussed in Example 2.15. 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 2.15 that the estimation of  $y$  is unaffected by a nonsingular transformation of the explanatory variables; that is, we

could just as well work with the model

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

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 utilizing 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  $W_1' W_1 = U' Z_1' Z_1 U = \Lambda$  and  $\mathbf{1}_N' W_1 = \mathbf{1}_N' Z_1 U = \mathbf{0}' U = \mathbf{0}'$ , the least squares estimate of  $\alpha_1$  is

$$\hat{\alpha}_1 = (W_1' W_1)^{-1} W_1' y = \Lambda^{-1} W_1' y,$$

while 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  $\lambda_i$ 's will be very small, while the variances of the corresponding  $\alpha_i$ '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 very little to the estimation of  $y$ . 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

$$y = \alpha_0 \mathbf{1}_N + W_{11} \alpha_{11} + \epsilon,$$

where  $W_{11}$  and  $\alpha_{11}$  are obtained from  $W_1$  and  $\alpha_1$  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  $\alpha_{11}$  can be written as

$$\hat{\alpha}_{11} = (W_{11}' W_{11})^{-1} W_{11}' y = \Lambda_1^{-1} W_{11}' y$$

Note that due to 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  $\delta_1$  in the original standardized model. Recall that  $\delta_1$  and  $\alpha_1$  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

$\delta_1 = U_1 \alpha_{11}$ , and so the principal components regression estimate of  $\delta_1$  is given by  $\hat{\delta}_{11} = U_1 \hat{\alpha}_{11} = U_1 \Lambda_1^{-1} W'_{11} y$ .

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  $\lambda$  be an eigenvalue of the  $m \times m$  symmetric matrix  $A$  with multiplicity  $r \geq 1$ . If  $x_1, \dots, x_r$  is a set of orthonormal eigenvectors corresponding to  $\lambda$ , then the eigenprojection of  $A$  associated with the eigenvalue  $\lambda$  is given by

$$P_A(\lambda) = \sum_{i=1}^r x_i x_i'$$

The eigenprojection  $P_A(\lambda)$  is simply the projection matrix for the vector space  $S_A(\lambda)$ . Thus, for any  $x \in R^m$ ,  $y = P_A(\lambda)x$  gives the orthogonal projection of  $x$  onto the eigenspace  $S_A(\lambda)$ . If we define  $X_1$  as before, that is  $X_1 = (x_1, \dots, x_r)$ , then  $P_A(\lambda) = X_1 X_1'$ . Note that  $P_A(\lambda)$  is unique even though the set of eigenvectors  $x_1, \dots, x_r$  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  $\lambda$ , but

$$Y_1 Y_1' = (X_1 Q)(X_1 Q)' = X_1 Q Q' X_1' = X_1 I_r X_1' = X_1 X_1' = P_A(\lambda)$$

The term spectral decomposition comes from the use of 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  $k \leq m$ , since some of the  $\mu_i$  may correspond to multiple eigenvalues. The set of  $\mu_i$  may be different from our set of  $\lambda_i$  in that we do not repeat values for the  $\mu_i$ . Thus, if  $A$  is  $4 \times 4$  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  $\Lambda$  as previously defined, the spectral decomposition states that

$$A = X \Lambda X' = \sum_{i=1}^m \lambda_i x_i x_i' = \sum_{i=1}^k \mu_i P_A(\mu_i),$$

so that  $A$  has been decomposed into a sum of terms, one corresponding to each value in the spectral set.

**Example 3.7.** It can be easily verified by solving the characteristic equation for the  $3 \times 3$  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  $(1/\sqrt{3}, 1/\sqrt{3}, 1/\sqrt{3})'$ , while 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

$$\begin{aligned} \begin{bmatrix} 5 & -1 & -1 \\ -1 & 5 & -1 \\ -1 & -1 & 5 \end{bmatrix} &= \begin{bmatrix} 1/\sqrt{3} & -2/\sqrt{6} & 0 \\ 1/\sqrt{3} & 1/\sqrt{6} & 1/\sqrt{2} \\ 1/\sqrt{3} & 1/\sqrt{6} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} 3 & 0 & 0 \\ 0 & 6 & 0 \\ 0 & 0 & 6 \end{bmatrix} \\ &\times \begin{bmatrix} 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ -2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix}, \end{aligned}$$

and the two eigenprojections of  $A$  are

$$\begin{aligned} P_A(3) &= \begin{bmatrix} 1/\sqrt{3} \\ 1/\sqrt{3} \\ 1/\sqrt{3} \end{bmatrix} [1/\sqrt{3} \quad 1/\sqrt{3} \quad 1/\sqrt{3}] = \frac{1}{3} \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \\ P_A(6) &= \begin{bmatrix} -2/\sqrt{6} & 0 \\ 1/\sqrt{6} & 1/\sqrt{2} \\ 1/\sqrt{6} & -1/\sqrt{2} \end{bmatrix} \begin{bmatrix} -2/\sqrt{6} & 1/\sqrt{6} & 1/\sqrt{6} \\ 0 & 1/\sqrt{2} & -1/\sqrt{2} \end{bmatrix} \\ &= \frac{1}{3} \begin{bmatrix} 2 & -1 & -1 \\ -1 & 2 & -1 \\ -1 & -1 & 2 \end{bmatrix} \end{aligned}$$

The relationship between the rank of a matrix and the number of its nonzero eigenvalues becomes an exact one for symmetric matrices.

**Theorem 3.11.** 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  $\Lambda$  has  $r$  nonzero diagonal elements and

$$\text{rank}(A) = \text{rank}(X\Lambda X') = \text{rank}(\Lambda),$$

since 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 3.8.** 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 in which we have equal variances and equal correlations; that is, the covariance matrix has the form

$$\Omega = \sigma^2 \begin{bmatrix} 1 & \rho & \cdots & \rho \\ \rho & 1 & \cdots & \rho \\ \vdots & \vdots & \ddots & \vdots \\ \rho & \rho & \cdots & 1 \end{bmatrix}$$

Alternatively,  $\Omega$  can be expressed as  $\Omega = \sigma^2 \{(1 - \rho)\mathbf{I}_m + \rho\mathbf{1}_m\mathbf{1}_m'\}$  so that it is a function of the vector  $\mathbf{1}_m$ . This vector also plays a crucial role in the eigenanalysis of  $\Omega$  since

$$\Omega \mathbf{1}_m = \sigma^2 \{(1 - \rho)\mathbf{1}_m + \rho\mathbf{1}_m\mathbf{1}_m'\mathbf{1}_m\} = \sigma^2 \{(1 - \rho) + m\rho\}\mathbf{1}_m$$

Thus,  $\mathbf{1}_m$  is an eigenvector of  $\Omega$  corresponding to the eigenvalue  $\sigma^2 \{(1 - \rho) + m\rho\}$ . The remaining eigenvalues of  $\Omega$  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  $\Omega$  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  $\rho$ ;  $\sigma^2 \{(1 - \rho) + m\rho\}$  will be larger than  $\sigma^2(1 - \rho)$  only if  $\rho$  is positive.

**Example 3.9.** 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, there is a much more limited set of possible eigenvalues and eigenvectors for correlation matrices. For the most extreme case, consider a  $2 \times 2$  correlation matrix that must have the form

$$P = \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$$

with  $-1 \leq \rho \leq 1$ , since  $P$  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  $Px = \lambda x$  we find that regardless of the value of  $\rho$ ,  $(1/\sqrt{2}, 1/\sqrt{2})'$  must be an eigenvector corresponding to  $1 + \rho$ , while  $(1/\sqrt{2}, -1/\sqrt{2})'$  must be an eigenvector corresponding to  $1 - \rho$ . Thus, ignoring sign changes, there is only one set of orthonormal eigenvectors possible for a  $2 \times 2$  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

$$P = \lambda_1 x_1 x_1' + \lambda_2 x_2 x_2' + \sum_{i=3}^m \lambda_i x_i x_i'$$

where  $\lambda_1$ ,  $\lambda_2$ , and  $\lambda$  are the distinct eigenvalues of  $P$ , and  $x_1, \dots, x_m$  are corresponding normalized eigenvectors. Since  $P$  is nonnegative definite, we must have  $\lambda_1 \geq 0$ ,  $\lambda_2 \geq 0$ , and  $\lambda \geq 0$ , while  $\text{tr}(P) = m$  implies that  $\lambda = (m - \lambda_1 - \lambda_2)/(m - 2)$ . Note that  $P$  can be written as

$$P = (\lambda_1 - \lambda)x_1 x_1' + (\lambda_2 - \lambda)x_2 x_2' + \lambda I_m,$$

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,

$$x_{i2}^2 = \frac{1 - \lambda - (\lambda_1 - \lambda)x_{i1}^2}{(\lambda_2 - \lambda)}$$

The constraints described can then be used to construct a particular matrix. For

instance, suppose that we want to construct a  $4 \times 4$  correlation matrix with eigenvalues  $\lambda_1 = 2$ ,  $\lambda_2 = 1$ , and  $\lambda = 0.5$  repeated twice. If we choose  $x_1 = (\frac{1}{2}, \frac{1}{2}, \frac{1}{2}, \frac{1}{2})'$ , then we must have  $x_2^2 = \frac{1}{4}$ , and so because of the orthogonality of  $x_1$  and  $x_2$ ,  $x_2$  can be any vector obtained from  $x_1$  by negating two of its components. For example, if we take  $x_2 = (\frac{1}{2}, -\frac{1}{2}, \frac{1}{2}, -\frac{1}{2})'$ , then

$$P = \begin{bmatrix} 1 & 0.25 & 0.50 & 0.25 \\ 0.25 & 1 & 0.25 & 0.50 \\ 0.50 & 0.25 & 1 & 0.25 \\ 0.25 & 0.50 & 0.25 & 1 \end{bmatrix}$$

## 5. CONTINUITY OF EIGENVALUES AND EIGENPROJECTIONS

Our first result of this section is one which 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 this theorem can be found in Ostrowski (1973). For some other similar bounds see Elsner (1982).

**Theorem 3.12.** Let  $A$  and  $B$  be  $m \times m$  matrices possessing 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 3.12 will allow us to establish a very 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 3.12. 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 3.12 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  $i$ , there exists some  $j$  such that  $\gamma_{j,n}$  is very close to  $\lambda_i$ , or more precisely, as  $B_n \rightarrow A$ , the eigenvalues of  $B_n$  are converging to those of  $A$ . This leads to the following important result.

**Theorem 3.13.** Let  $\lambda_1, \dots, \lambda_m$  be the eigenvalues of the  $m \times m$  matrix  $A$ . Then, for each  $i$ ,  $\lambda_i$  is a continuous function of the elements of  $A$ .

Our next result 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 3.14.** Suppose that  $A$  is an  $m \times m$  symmetric matrix and  $\lambda$  is one of its eigenvalues. Then  $P_A(\lambda)$ , the eigenprojection associated with the eigenvalue  $\lambda$ , is a continuous function of the elements of  $A$ .

**Example 3.10.** Consider the matrix  $A$

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

which clearly has the simple eigenvalue 2 and the repeated eigenvalue 1. Suppose that  $B_1, B_2, \dots$  is a sequence of  $3 \times 3$  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  $B_n$ , while  $\mathbf{x}_{1,n}$ ,  $\mathbf{x}_{2,n}$ , and  $\mathbf{x}_{3,n}$  is a corresponding set of orthonormal eigenvectors. Theorem 3.13 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 3.14 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

$$1, \frac{3}{2} - \frac{\sqrt{1 + 4n^{-2}}}{2}, \quad \frac{3}{2} + \frac{\sqrt{1 + 4n^{-2}}}{2},$$

which do converge to 1, 1, and 2, respectively. It is left as an exercise for the reader to verify that

$$P_{1,n} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} = P_A(2), \quad P_{2,n} \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} = P_A(1)$$

## 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  $x'Ax$  as a function of  $x$ . If  $\alpha$  is a nonzero scalar, then  $(\alpha x)'A(\alpha x) = \alpha^2 x'Ax$ , so that the quadratic form can be made arbitrarily small or large, depending on whether  $x'Ax$  is negative or positive, through the proper choice of  $\alpha$ . Thus, any meaningful study of the variational properties of  $x'Ax$  as we change  $x$  will require the removal of the effect of scale changes in  $x$ . One way of doing this is through the construction of what is commonly called the Rayleigh quotient given by

$$R(x, A) = \frac{x'Ax}{x'x}$$

Note that  $R(\alpha x, A) = R(x, A)$ . Our first result involves the global maximization and minimization of  $R(x, A)$ .

**Theorem 3.15.** 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  $x \neq 0$ ,

$$\lambda_m \leq \frac{x'Ax}{x'x} \leq \lambda_1, \quad (3.6)$$

and, in particular,

$$\lambda_m = \min_{x \neq 0} \frac{x'Ax}{x'x}, \quad \lambda_1 = \max_{x \neq 0} \frac{x'Ax}{x'x} \quad (3.7)$$

*Proof.* Let  $A = X\Lambda X'$  be the spectral decomposition of  $A$ , where the columns of  $X = (x_1, \dots, x_m)$  are normalized eigenvectors of  $A$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ . Then, if  $y = X'x$ , we have

$$\frac{x'Ax}{x'x} = \frac{x'X\Lambda X'x}{x'XX'x} = \frac{y'\Lambda y}{y'y} = \frac{\sum_{i=1}^m \lambda_i y_i^2}{\sum_{i=1}^m y_i^2},$$

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  $x$  for which the bounds in (3.6) are attained; for instance, the lower bound is attained with  $x = x_m$ , while the upper bound holds with  $x = x_1$ .  $\square$

Note that, since for any nonnull  $x$ ,  $z = (x'x)^{-1/2}x$  is a unit vector, the minimization and maximization of  $z'Az$  over all unit vectors  $z$  will also yield  $\lambda_m$  and  $\lambda_1$ , respectively; that is,

$$\lambda_m = \min_{z'z=1} z'Az, \quad \lambda_1 = \max_{z'z=1} z'Az$$

The following theorem shows that each eigenvalue of a symmetric matrix  $A$  can be expressed as a constrained maximum or minimum of the Rayleigh quotient,  $R(x, A)$ .

**Theorem 3.16.** Let  $A$  be an  $m \times m$  symmetric matrix having eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  with  $x_1, \dots, x_m$  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  $X_h = (x_1, \dots, x_h)$  and  $Y_h = (x_h, \dots, x_m)$ , respectively. Then

$$\lambda_h = \min_{x \in S_h} \frac{x'Ax}{x'x} = \min_{Y_{h+1}'x=0} \frac{x'Ax}{x'x},$$

and

$$\lambda_h = \max_{x \in T_h} \frac{x'Ax}{x'x} = \max_{X_{h-1}'x=0} \frac{x'Ax}{x'x},$$

where the vector  $x = 0$  has been excluded from the maximization and minimization processes.

*Proof.* We will prove the result concerning the minimum; the proof for the maximum is similar. Let  $X = (x_1, \dots, x_m)$  and  $\Lambda = \text{diag}(\lambda_1, \dots, \lambda_m)$ . Note that, since  $X'AX = \Lambda$  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  $x \in S_h$  if and only if there exists an  $h \times 1$  vector  $y$  such that  $x = X_h y$ . Consequently,

$$\min_{x \in S_h} \frac{x'Ax}{x'x} = \min_{y \neq 0} \frac{y'X_h'AX_h y}{y'X_h'X_h y} = \min_{y \neq 0} \frac{y'\Lambda_h y}{y'y} = \lambda_h,$$

where the last equality follows from Theorem 3.15. The second version of the minimization follows immediately from the first and the fact that the null space of  $Y_{h+1}'$  is  $S_h$ .  $\square$

The next two examples give some indication of how the extremal properties of eigenvalues make them important features in many applications.

**Example 3.11.** 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  $\mu_1, \dots, \mu_k$  represent the  $k$  group mean vectors, and let  $\mu = (\mu_1 + \dots + \mu_k)/k$  be the average of these mean vectors. To investigate the differences in group means, we will utilize the deviations  $(\mu_i - \mu)$  from the average mean; in particular, we form the sum of squares and cross products matrix given by

$$A = \sum_{i=1}^k (\mu_i - \mu)(\mu_i - \mu)'$$

Note that for a particular unit vector  $x$ ,  $x'Ax$  will give a measure of the dif-

ferences among the  $k$  groups in the direction  $\mathbf{x}$ ; a value of zero indicates the groups have identical means in this direction, while 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 Theorems 3.15 and 3.16 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  $\lambda_1$  and  $\lambda_2$ . 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 4.11 we will discuss the statistical procedure, called canonical variate analysis, that utilizes this sort of dimension reducing process.

**Example 3.12.** In Example 3.11, the focus was on means. In this example, we will look at a procedure that concentrates on variances. This technique, called principal component analysis, was developed by Hotelling (1933). Some good references on this subject are Jackson (1991) and Jolliffe (1986). Let  $\mathbf{x}$  be an  $m \times 1$  random vector having the covariance matrix  $\Omega$ . 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. But 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 arbitrarily large by taking  $\mathbf{a}_1 = \alpha \mathbf{c}$  for some scalar  $\alpha$  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 3.15 that this direction is given by the normalized eigenvector of  $\Omega$  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 3.16, this second direction is given by the normalized eigenvector of  $\Omega$  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  $\Omega$ . 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  $\Omega$ , are uncorrelated since for  $i \neq j$ ,

$$\text{cov}(\mathbf{a}'_i \mathbf{x}, \mathbf{a}'_j \mathbf{x}) = \mathbf{a}'_i \Omega \mathbf{a}_j = 0$$

For some specific examples, first consider the  $4 \times 4$  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  $\Omega$  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  $\mathbf{x}$ . This 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  $\Omega$ ,  $(0.25, 0.25, 0.25, 0.25)'$  and  $(0.25, 0.25, -0.25, -0.25)'$ . 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  $\Omega$  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.

Our final result, 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 3.17.** 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

$$\lambda_h = \min_{B_h} \max_{B'_h \mathbf{x} = \mathbf{0}} \frac{\mathbf{x}' A \mathbf{x}}{\mathbf{x}' \mathbf{x}}, \quad (3.9)$$

as well as



$$\lambda_h = \max_{C_h} \min_{C_h'x=0} \frac{x'Ax}{x'x} \tag{3.10}$$

where the vector  $x = 0$  has been excluded.

*Proof.* We first prove the min-max result given by (3.9). Let  $X_h = (x_1, \dots, x_h)$ , where  $x_1, \dots, x_h$  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

$$\min_{B_h} \max_{B_h'x=0} \frac{x'Ax}{x'x} \leq \max_{X_{h-1}'x=0} \frac{x'Ax}{x'x} = \lambda_h, \tag{3.11}$$

where the equality follows from Theorem 3.16. 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  $y$  such that  $B_h'X_h y = 0$ . Since  $X_h y$  is one choice for  $x$ , we find that

$$\max_{B_h'x=0} \frac{x'Ax}{x'x} \geq \frac{y'X_h'AX_h y}{y'X_h'X_h y} = \frac{y'\Lambda_h y}{y'y} \geq \lambda_h, \tag{3.12}$$

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

$$\min_{B_h} \max_{B_h'x=0} \frac{x'Ax}{x'x} \geq \lambda_h$$

This, along with (3.11), proves (3.9). The proof of (3.10) is along the same lines. Let  $Y_h = (x_{h+1}, \dots, x_m)$ , where  $x_{h+1}, \dots, x_m$  is a set of orthonormal eigenvectors of  $A$ , corresponding to the eigenvalues  $\lambda_{h+1}, \dots, \lambda_m$ . Since  $Y_{h+1}$  is an  $m \times (m-h)$  matrix satisfying  $Y_{h+1}'Y_{h+1} = I_{m-h}$ , it follows that

$$\max_{C_h} \min_{C_h'x=0} \frac{x'Ax}{x'x} \geq \min_{Y_{h+1}'x=0} \frac{x'Ax}{x'x} = \lambda_h, \tag{3.13}$$

where the equality follows from Theorem 3.16. 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, there exists an  $(m-h+1) \times 1$  nonnull vector

$y$  satisfying  $C'_h Y_h y = 0$ . Since  $Y_h y$  is one choice for  $x$ , we have

$$\min_{C'_h x = 0} \frac{x'Ax}{x'x} \leq \frac{y'Y'_h A Y_h y}{y'Y'_h Y_h y} = \frac{y'\Delta_h y}{y'y} \leq \lambda_h, \quad (3.14)$$

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

$$\max_{C_h} \min_{C'_h x = 0} \frac{x'Ax}{x'x} \leq \lambda_h$$

This together with (3.13) establishes (3.10).  $\square$

**Corollary 3.17.1.** 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

$$\lambda_h \leq \max_{B'_h x = 0} \frac{x'Ax}{x'x},$$

and

$$\lambda_h \geq \min_{C'_h x = 0} \frac{x'Ax}{x'x}$$

*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 3.17. We need to establish them for arbitrary  $B_h$  and  $C_h$ . When  $B'_h B_h = I_{h-1}$ , the set  $S_{B_h} = \{x: x \in R^m, B'_h x = 0\}$  is the orthogonal complement of the vector space which has the columns of  $B_h$  as an orthonormal basis. Thus, the first inequality holds when maximizing over all  $x \neq 0$  in any  $(m-h+1)$ -dimensional vector subspace of  $R^m$ . Consequently, this inequality also will hold for any  $m \times (h-1)$  matrix  $B_h$  since, 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 result is left to the reader as an exercise.

**Theorem 3.18.** Suppose that  $A$  and  $B$  are  $m \times m$  symmetric matrices and  $A - B$  is nonnegative definite. Then  $\lambda_i(A) \geq \lambda_i(B)$  for  $i = 1, \dots, m$ .

Some additional extremal properties of eigenvalues can be found in Bellman (1970) and Horn and Johnson (1985).

## 7. SOME ADDITIONAL RESULTS CONCERNING EIGENVALUES

Let  $A$  be an  $m \times m$  symmetric matrix and  $H$ , an  $m \times h$  matrix satisfying  $H'H = I_h$ . In some situations it is of interest to compare the eigenvalues of  $A$  to those of  $H'AH$ . Some comparisons follow immediately from Theorem 3.17. 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

$$\lambda_h(H'AH) \leq \lambda_h(A)$$

The following result, 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

$$\begin{aligned} \lambda_{m-h+i}(A) &= \lambda_{n-1}(A) = \min_{Y_n'x=0} \frac{\mathbf{x}'Ax}{\mathbf{x}'x} \leq \min_{\substack{Y_n'x=0 \\ x=Hy}} \frac{\mathbf{x}'Ax}{\mathbf{x}'x} \\ &= \min_{Y_n'Hy=0} \frac{\mathbf{y}'H'AH\mathbf{y}}{\mathbf{y}'\mathbf{y}} \leq \lambda_{h-(m-n+1)}(H'AH) = \lambda_i(H'AH), \end{aligned}$$

where the second equality follows from Theorem 3.16. The last inequality follows from Corollary 3.17.1, 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

$$\begin{aligned}\lambda_i(A) &= \max_{x'_{i-1}x=0} \frac{x'Ax}{x'x} \geq \max_{\substack{x'_{i-1}x=0 \\ x=Hy}} \frac{x'Ax}{x'x} \\ &= \max_{x'_{i-1}Hy=0} \frac{y'H'AHy}{y'y} \geq \lambda_i(H'AH),\end{aligned}$$

where the first equality follows from Theorem 3.16 and the final inequality follows from Corollary 3.17.1.  $\square$

Theorem 3.19 can be used to prove the following useful result.

**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  $i = 1, \dots, k$ ,

$$\lambda_{m-i+1}(A) \leq \lambda_{k-i+1}(A_k) \leq \lambda_{k-i+1}(A)$$

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  $x'Ax$ . We now show that these conditions also can be expressed in terms of the eigenvalues of  $A$ .

**Theorem 3.21.** Let  $A$  be an  $m \times m$  symmetric matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$ . 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  $X = (x_1, \dots, x_m)$  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  $x'Ax > 0$  for all  $x \neq 0$ , so in particular, choosing  $x = x_i$ , we have

$$x'_i A x_i = \lambda_i > 0$$

Conversely, if  $\lambda_i > 0$  for all  $i$ , then for any  $x \neq 0$  define  $y = X'x$ , and note that

$$x'Ax = x'X\Lambda X'x = y'\Lambda y = \sum_{i=1}^m y_i^2 \lambda_i \tag{3.15}$$

has to be positive because the  $\lambda_i$  are positive and at least one of the  $y_i^2$  is pos-

itive since  $\mathbf{y} \neq \mathbf{0}$ . 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 (3.15) 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 3.21 that positive definite matrices are nonsingular, while positive semidefinite matrices are singular.

**Example 3.13.** Consider the ordinary least squares estimator  $\hat{\boldsymbol{\beta}} = (X'X)^{-1}X'y$  of  $\boldsymbol{\beta}$  in the model

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

where  $E(\boldsymbol{\epsilon}) = \mathbf{0}$  and  $\text{var}(\boldsymbol{\epsilon}) = \sigma^2 I_N$ . For an arbitrary  $(k+1) \times 1$  vector  $\mathbf{c}$ , we will prove that  $\mathbf{c}'\hat{\boldsymbol{\beta}}$  is the best linear unbiased estimator of  $\mathbf{c}'\boldsymbol{\beta}$ ; an estimator  $t$  is an unbiased estimator of  $\mathbf{c}'\boldsymbol{\beta}$  if  $E(t) = \mathbf{c}'\boldsymbol{\beta}$ . Clearly,  $\mathbf{c}'\hat{\boldsymbol{\beta}}$  is unbiased since  $E(\boldsymbol{\epsilon}) = \mathbf{0}$  implies that

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

To show that it 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  $\mathbf{c}'\boldsymbol{\beta}$ . Let  $\mathbf{a}'\mathbf{y}$  be an arbitrary linear unbiased estimator of  $\mathbf{c}'\boldsymbol{\beta}$ , so that

$$\mathbf{c}'\boldsymbol{\beta} = E(\mathbf{a}'\mathbf{y}) = \mathbf{a}'E(\mathbf{y}) = \mathbf{a}'X\boldsymbol{\beta},$$

regardless of the value of the vector  $\boldsymbol{\beta}$ . But this implies that

$$\mathbf{c}' = \mathbf{a}'X$$

Now

$$\text{var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) = \mathbf{c}'\{\text{var}(\hat{\boldsymbol{\beta}})\}\mathbf{c} = \mathbf{c}'\{\sigma^2(X'X)^{-1}\}\mathbf{c} = \sigma^2\mathbf{a}'X(X'X)^{-1}X'\mathbf{a},$$

while

$$\text{var}(\mathbf{a}'\mathbf{y}) = \mathbf{a}'\{\text{var}(\mathbf{y})\}\mathbf{a} = \mathbf{a}'\{\sigma^2 I_N\}\mathbf{a} = \sigma^2\mathbf{a}'\mathbf{a}$$

Thus, the difference in their variances is

$$\begin{aligned}\text{var}(\mathbf{a}'\mathbf{y}) - \text{var}(\mathbf{c}'\hat{\boldsymbol{\beta}}) &= \sigma^2\mathbf{a}'\mathbf{a} - \sigma^2\mathbf{a}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{a} \\ &= \sigma^2\mathbf{a}'(\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')\mathbf{a}\end{aligned}$$

But

$$\{\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\}^2 = \{\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\},$$

and so using Theorem 3.4, we find that each of the eigenvalues of  $\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  must be 0 or 1. Thus, from Theorem 3.21, we see that  $\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$  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 3.22.** 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 3.21 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 3.23.** 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  $\mathbf{X}$ , 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  $\lambda$  is also an eigenvalue of  $TT'$ . Its multiplicity is also  $h$  since

$$\begin{aligned} \text{rank}(Y) &= \text{rank}(TX) = \text{rank}((TX)'TX) = \text{rank}(X'T'TX) \\ &= \text{rank}(\lambda X'X) = \text{rank}(\lambda I_h) = h \end{aligned} \quad \square$$

Next we will use the Courant–Fischer min-max theorem to prove the following important monotonicity property of the eigenvalues of symmetric matrices.

**Theorem 3.24.** 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

$$\begin{aligned} \max_{B_h'x=0} \frac{x'(A+B)x}{x'x} &= \max_{B_h'x=0} \left( \frac{x'Ax}{x'x} + \frac{x'Bx}{x'x} \right) \geq \max_{B_h'x=0} \frac{x'Ax}{x'x} \\ &\quad + \min_{B_h'x=0} \frac{x'Bx}{x'x} \\ &\geq \max_{B_h'x=0} \frac{x'Ax}{x'x} + \min_{x \neq 0} \frac{x'Bx}{x'x} \\ &= \max_{B_h'x=0} \frac{x'Ax}{x'x} + \lambda_m(B) \geq \max_{B_h'x=0} \frac{x'Ax}{x'x}, \end{aligned}$$

where the last equality follows from Theorem 3.15. The final inequality above is strict if  $B$  is positive definite since, 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 3.17, we get

$$\lambda_h(A + B) = \min_{B_h} \max_{B_h'x=0} \frac{x'(A + B)x}{x'x} \geq \min_{B_h} \max_{B_h'x=0} \frac{x'Ax}{x'x} = \lambda_h(A)$$

This completes the proof. □

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

$$\lambda_2(A + B) = 8 < \lambda_2(A) + \lambda_2(B) = 3 + 6 = 9,$$

while

$$\lambda_3(A + B) = 7 > \lambda_3(A) + \lambda_3(B) = 2 + 4 = 6$$

In Example 3.11 we discussed a situation in which the eigenvalues and eigenvectors of

$$A = \sum_{i=1}^k (\mu_i - \mu)(\mu_i - \mu)'$$

were utilized in analyzing differences among the group means  $\mu_1, \dots, \mu_k$ . For instance, an eigenvector  $x_1$ , corresponding to the largest eigenvalue of  $A$ , gives the direction of maximum dispersion among the group means in that

$$\frac{x_1'Ax_1}{x_1'x_1}$$

is maximized. The division here by  $x_1'x_1$ , 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  $y$  is a random vector with covariance matrix  $B$ , then the variability of  $y$  in the direction given by  $x$  will be  $\text{var}(x'y) = x'Bx$ . 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

$$\frac{x'Ax}{x'Bx}$$

The vector  $x_1$  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  $x_1$  would be to find the vector  $x_2$  that maximizes the ratio but has  $x_2'y$  uncorrelated with  $x_1'y$ ; this would be the vector  $x_2$  that maximizes the ratio in the equation subject to the constraint that  $x_1'Bx_2 = 0$ . Continuing in this fashion, we would determine the  $m$  vectors  $x_1, \dots, x_m$  that yield the  $m$  extremal values  $\lambda_1, \dots, \lambda_m$  of the ratio. These extremal values are identified in the following theorem.

**Theorem 3.25.** Let  $A$  and  $B$  be  $m \times m$  matrices, with  $A$  being nonnegative definite and  $B$  positive definite. For  $h = 1, \dots, m$ , define  $X_h = (x_1, \dots, x_h)$  and  $Y_h = (x_h, \dots, x_m)$ , where  $x_1, \dots, x_m$  are linearly independent eigenvectors of  $B^{-1}A$  corresponding to the eigenvalues  $\lambda_1(B^{-1}A) \geq \dots \geq \lambda_m(B^{-1}A)$ . Then

$$\lambda_h(B^{-1}A) = \min_{Y'_{h+1}Bx=0} \frac{x'Ax}{x'Bx},$$

and

$$\lambda_h(B^{-1}A) = \max_{X'_{h-1}Bx=0} \frac{x'Ax}{x'Bx},$$

where  $x = 0$  is excluded, and the min and max are over all  $x \neq 0$  when  $h = m$  and  $h = 1$ , respectively.

*Proof.* We will prove the result involving the minimum; the proof for the maximum is similar. 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 due to Theorem 3.19. 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. Putting  $y = Tx$ , we find that

$$\begin{aligned} \min_{Y'_{h+1}Bx=0} \frac{x'Ax}{x'Bx} &= \min_{Y'_{h+1}TTx=0} \frac{x'TT^{-1}AT^{-1}Tx}{x'TTx} \\ &= \min_{Y'_{h+1}Ty=0} \frac{y'T^{-1}AT^{-1}y}{y'y} \end{aligned} \tag{3.16}$$

Note that if we write  $\lambda_i = \lambda_i(B^{-1}A)$ , then  $B^{-1}Ax_i = \lambda_i x_i$ , so that

$$T^{-1}T^{-1}Ax_i = \lambda_i x_i,$$

which implies

$$T^{-1}AT^{-1}Tx_i = \lambda_i Tx_i$$

Thus,  $Tx_i$  is an eigenvector of  $T^{-1}AT^{-1}$  corresponding to the eigenvalue

$\lambda_i = \lambda_i(T^{-1}AT^{-1})$ ; that is, the eigenvalues of  $B^{-1}A$  are the same as those of  $T^{-1}AT^{-1}$ . Since the rows of  $Y'_{h+1}T$  are the transposes of the eigenvectors  $Tx_{h+1}, \dots, Tx_m$ , it follows from Theorem 3.16 that (3.16) equals  $\lambda_h(T^{-1}AT^{-1})$ , which we have already established as being the same as  $\lambda_h(B^{-1}A)$ .  $\square$

Since  $x_i$  is an eigenvector of  $B^{-1}A$  corresponding to the eigenvalue  $\lambda_i = \lambda_i(B^{-1}A)$ , we know that

$$B^{-1}Ax_i = \lambda_i x_i$$

or, equivalently,

$$Ax_i = \lambda_i Bx_i \quad (3.17)$$

Equation (3.17) is similar to the eigenvalue–eigenvector equation of  $A$ , except for the multiplication of  $x_i$  by  $B$  on the right-hand side of the equation. The eigenvalues satisfying (3.17) are sometimes referred to as the eigenvalues of  $A$  in the metric of  $B$ . Note that if we premultiply (3.17) by  $x'_i$  and then solve for  $\lambda_i$ , we get

$$\lambda_i(B^{-1}A) = \frac{x'_i Ax_i}{x'_i Bx_i};$$

that is, the extremal values given in Theorem 3.25 are attained at the eigenvectors of  $B^{-1}A$ .

The proof of the previous theorem suggests a way of simultaneously diagonalizing the matrices  $A$  and  $B$ . Since  $T^{-1}AT^{-1}$  is symmetric, it can be expressed in the form  $Q\Lambda Q'$ , where  $Q$  is orthogonal and  $\Lambda = \text{diag}(\lambda_1(T^{-1}AT^{-1}), \dots, \lambda_m(T^{-1}AT^{-1}))$ . The matrix  $C = Q'T^{-1}$  is nonsingular since  $Q$  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. For some other related results see Olkin and Tomsky (1981).

**Theorem 3.26.** Let  $A$  be an  $m \times m$  nonnegative definite 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  $i = 1, \dots, h$

$$\lambda_i((F'AF)(F'BF)^{-1}) \leq \lambda_i(AB^{-1}),$$

and further

$$\max_F \lambda_i((F'AF)(F'BF)^{-1}) = \lambda_i(AB^{-1})$$

*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

$$\begin{aligned} \max_F \lambda_i((F'AF)(F'BF)^{-1}) &= \max_F \lambda_i((F'G\Lambda G'F)(F'GG'F)^{-1}) \\ &= \max_E \lambda_i((E'\Lambda E)(E'E)^{-1}), \end{aligned}$$

where this last maximization is also over all  $m \times h$  matrices of rank  $h$ , since  $E = G'F$  must have the same rank as  $F$ . Note that since  $E$  has rank  $h$ , the  $h \times h$  matrix  $E'E$  is a nonsingular symmetric matrix. As was seen in the previous proof, such a matrix can be expressed as  $E'E = TT$  for some nonsingular symmetric  $h \times h$  matrix  $T$ . It then follows that

$$\begin{aligned} \max_E \lambda_i((E'\Lambda E)(E'E)^{-1}) &= \max_E \lambda_i((E'\Lambda E)(TT)^{-1}) \\ &= \max_E \lambda_i(T^{-1}E'\Lambda ET^{-1}), \end{aligned}$$

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  $H'H = T^{-1}E'ET^{-1} = T^{-1}TTT^{-1} = I_h$ . 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 3.14.** 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  $\mu_i$ s are constants and the  $\epsilon_{ij}$ s are independent and identically distributed as  $N(0, \sigma^2)$ . Our goal is to determine whether or not the  $\mu_i$ 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  $\mu_i$ s differ. An analysis of variance compares (see Problem 2.31) the variability between treatments,

$$SST = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})^2,$$

to the variability within treatments,

$$SSE = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2,$$

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 while SSE is called the sum of squares for error. The hypothesis  $H_0$  is rejected if the statistic

$$F = \frac{SST/(k-1)}{SSE/(n-k)}$$

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  $y_{ij}$  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

$$y_{ij} = \mu_i + \epsilon_{ij},$$

where  $\mu_i$  is an  $m \times 1$  vector of constants and  $\epsilon_{ij} \sim N_m(\mathbf{0}, \Omega)$ , independently. Measures of the between treatment variability and within treatment variability are now given by the matrices,

$$B = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})(\bar{y}_i - \bar{y})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)(y_{ij} - \bar{y}_i)'$$

One approach to testing the null hypothesis  $H_0: \mu_1 = \dots = \mu_k$  against the alternative hypothesis,  $H_1$ : at least two of the  $\mu_i$ 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  $c$  is any  $m \times 1$  vector, and we define the hypothesis  $H_0(c): c' \mu_1 = \dots = c' \mu_k$ , then the intersection of  $H_0(c)$  over all  $c \in R^m$  is the hypothesis  $H_0$ . In addition, if we define the hypothesis  $H_1(c)$ : at least two of the  $c' \mu_i$ s differ, then the union of the hypotheses  $H_1(c)$  over all  $c \in R^m$  is the hypothesis  $H_1$ . Thus, we should reject the hypothesis  $H_0$  if and only if we reject  $H_0(c)$  for at least one  $c$ . Now the null hypothesis  $H_0(c)$  involves the univariate one-way classification model in which  $c' y_{ij}$  is the response, and so we would reject  $H_0(c)$  for large values of the  $F$  statistic

$$F(c) = \frac{SST(c)/(k-1)}{SSE(c)/(n-k)},$$

where  $SST(c)$  and  $SSE(c)$  are the sums of squares for treatments and errors, respectively, computed for the responses  $c' y_{ij}$ . Since  $H_0$  is rejected if  $H_0(c)$  is rejected for at least one  $c$ , we will reject  $H_0$  if  $F(c)$  is sufficiently large for at least one  $c$  or, equivalently, if

$$\max_{c \neq 0} F(c)$$

is sufficiently large. Omitting the constants  $(k-1)$  and  $(n-k)$  and noting that the sums of squares  $SST(c)$  and  $SSE(c)$  can be expressed using  $B$  and  $W$  as

$$SST(c) = c' Bc, \quad SSE(c) = c' Wc,$$

we find that we reject  $H_0$  for large values of

$$\max_{c \neq 0} \frac{c' Bc}{c' Wc} = \lambda_1(W^{-1}B), \quad (3.18)$$

where the right-hand side follows from Theorem 3.25. 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 (1990)] so that

$$P[\lambda_1(W^{-1}B) \leq u_{1-\alpha} | H_0] = 1 - \alpha, \quad (3.19)$$

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 (3.18) and (3.19) that for any mean

vectors  $\mu_1, \dots, \mu_k$ , with probability  $1 - \alpha$ , the inequality

$$\frac{\sum_{i=1}^k n_i \mathbf{c}' \{(\bar{y}_i - \bar{y}) - (\mu_i - \mu)\} \{(\bar{y}_i - \bar{y}) - (\mu_i - \mu)\}' \mathbf{c}}{\mathbf{c}' W \mathbf{c}} \leq u_{1-\alpha}, \quad (3.20)$$

holds for all  $m \times 1$  vectors  $\mathbf{c}$ , where

$$\mu = \sum_{i=1}^k n_i \mu_i / n$$

Scheffé's method [see Scheffé (1953) or Miller (1981)] can then be used on (3.20) to yield the inequalities

$$\begin{aligned} \sum_{i=1}^k \sum_{j=1}^m a_i c_j \bar{x}_{ij} - \sqrt{u_{1-\alpha} \mathbf{c}' W \mathbf{c} \left( \sum_{i=1}^k a_i^2 / n_i \right)} \\ \leq \sum_{i=1}^k \sum_{j=1}^m a_i c_j \mu_{ij} \\ \leq \sum_{i=1}^k \sum_{j=1}^m a_i c_j \bar{x}_{ij} + \sqrt{u_{1-\alpha} \mathbf{c}' W \mathbf{c} \left( \sum_{i=1}^k a_i^2 / n_i \right)} \end{aligned}$$

which hold with probability  $1 - \alpha$ , for all  $m \times 1$  vectors  $\mathbf{c}$  and all  $k \times 1$  vectors  $\mathbf{a}$  satisfying  $\mathbf{a}' \mathbf{1}_k = 0$ .

## PROBLEMS

1. Consider the  $3 \times 3$  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.
2. Find the eigenvalues of  $A'$ , where  $A$  is the matrix given in Problem 1. Determine the eigenspaces for  $A'$  and compare these to those of  $A$ .

3. Let the  $3 \times 3$  matrix  $A$  be given by

$$A = \begin{bmatrix} 1 & -2 & 0 \\ 1 & 4 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

- Find the eigenvalues of  $A$ .
  - For each different value of  $\lambda$ , determine the associated eigenspace  $S_A(\lambda)$ .
  - Describe the eigenspaces obtained in part (b).
4. 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  $(A + \gamma I)$  has eigenvalues  $\lambda_1 + \gamma, \dots, \lambda_m + \gamma$  and corresponding eigenvectors  $\mathbf{x}_1, \dots, \mathbf{x}_m$ .
5. In Example 3.6, 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  $\hat{\delta}_1 = (Z_1' Z_1)^{-1} Z_1' y$  by  $\hat{\delta}_{1\gamma} = (Z_1' Z_1 + \gamma I)^{-1} Z_1' y$ , where  $\gamma$  is a small positive number. This adjustment will reduce the impact of the near singularity of  $Z_1' Z_1$  since the addition of  $\gamma I$  increases each of the eigenvalues of  $Z_1' Z_1$  by  $\gamma$ .
- Show that if  $N > 2k + 1$ , there is an  $N \times k$  matrix  $W$  such that  $\hat{\delta}_{1\gamma}$  is the ordinary least squares estimate of  $\delta_1$  in the model

$$y = \delta_0 \mathbf{1}_N + (Z_1 + W) \delta_1 + \epsilon;$$

that is,  $\hat{\delta}_{1\gamma}$  can be viewed as the ordinary least squares estimator of  $\delta_1$  after we have perturbed the matrix of values for the explanatory variables  $Z_1$  by  $W$ .

- Show that there exists a  $k \times k$  matrix  $U$  such that  $\hat{\delta}_{1\gamma}$  is the ordinary least squares estimate of  $\delta_1$  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  $\epsilon_* \sim N_k(\mathbf{0}, \sigma^2 I)$ , independently of  $\epsilon$ . 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.

6. Refer to Example 3.6 and the previous exercise.
- 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  $\delta_1$ .
  - 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  $\delta_1$ .
  - 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.
7. 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.
8. If  $\lambda$  is a real eigenvalue of the  $m \times m$  real matrix  $A$ , show that there exist real eigenvectors of  $A$  corresponding to the eigenvalue  $\lambda$ .
9. Prove the results given in Theorem 3.2.
10. Suppose that  $\lambda$  is a simple eigenvalue of the  $m \times m$  matrix  $A$ . Show that  $\text{rank}(A - \lambda I) = m - 1$ .
11. If  $A$  is an  $m \times m$  matrix and  $\text{rank}(A - \lambda I) = m - 1$ , show that  $\lambda$  is an eigenvalue of  $A$  with multiplicity of at least one.
12. Consider the  $m \times m$  matrix

$$A = \begin{bmatrix} 1 & 1 & 0 & \cdots & 0 \\ 0 & 1 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 0 & 0 & 0 & \cdots & 1 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix},$$

which has each element on and directly above the diagonal equal to 1. Find the eigenvalues and eigenvectors of  $A$ .

13. 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  $I + A$  is nonsingular, find the eigenvalues and eigenvectors of
- $(I + A)^{-1}$ ,
  - $A + A^{-1}$ ,
  - $(I + A^{-1})$ .



## PROBLEMS

14. Let the  $m \times m$  nonsingular matrix  $A$  be such that  $I + A$  is nonsingular, and define

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

- (a) Show that if  $x$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\lambda$ , then  $x$  is an eigenvector of  $B$  corresponding to the eigenvalue 1.  
 (b) Use Theorem 1.7 to show that  $B = I$ .

15. Consider the  $2 \times 2$  matrix

$$A = \begin{bmatrix} 4 & 2 \\ 3 & 5 \end{bmatrix}$$

- (a) Find the characteristic equation of  $A$ .  
 (b) Illustrate Theorem 3.7 by substituting  $A$  for  $\lambda$  in the characteristic equation obtained in (a) and then showing that 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  $I$ .  
 (d) In a similar fashion, write  $A^3$  and  $A^{-1}$  as linear combinations of  $A$  and  $I$ .

16. Consider the general  $2 \times 2$  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?

17. Find the eigenvalues and eigenvectors of the matrix  $\mathbf{1}_m \mathbf{1}'_m$ .

18. Consider the  $m \times m$  matrix  $A = \alpha \mathbf{I}_m + \beta \mathbf{1}_m \mathbf{1}'_m$ , where  $\alpha$  and  $\beta$  are scalars.  
 (a) Find the eigenvalues and eigenvectors of  $A$ .  
 (b) Determine the eigenspaces and associated eigenprojections of  $A$ .  
 (c) For which values of  $\alpha$  and  $\beta$  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  $\alpha^{m-1}(\alpha + m\beta)$ .

19. Consider the  $m \times m$  matrix  $A = \alpha I_m + \beta c c'$ , where  $\alpha$  and  $\beta$  are scalars and  $c \neq \mathbf{0}$  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  $A$ .

20. Let  $A$  be the  $3 \times 3$  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  $A$ .

21. Construct a  $3 \times 3$  symmetric matrix having eigenvalues 18, 21, and 28, and corresponding eigenvectors  $(1, 1, 2)'$ ,  $(4, -2, -1)'$ , and  $(1, 3, -2)'$ .

22. Show that if  $A$  is an  $m \times m$  symmetric matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$ , then

$$\sum_{i=1}^m \sum_{j=1}^m a_{ij}^2 = \sum_{i=1}^m \lambda_i^2$$

23. 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.

24. Use Theorem 3.17 to prove Theorem 3.18. Show that the converse is not true; that is, find symmetric matrices  $A$  and  $B$  for which  $\lambda_i(A) \geq \lambda_i(B)$  for  $i = 1, \dots, m$ , yet  $A - B$  is not nonnegative definite.

25. Let  $A$  be an  $m \times n$  matrix with  $\text{rank}(A) = r$ . Use the spectral decomposition of  $A'A$  to show that there exists an  $n \times (n - r)$  matrix  $X$  such that

$$AX = (0) \quad \text{and} \quad X'X = I_{n-r}$$

In a similar fashion, show that there exists an  $(m-r) \times m$  matrix  $Y$  such that

$$YA = (0) \quad \text{and} \quad YY' = I_{m-r}$$

26. Let  $A$  be the  $2 \times 3$  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.

27. 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.

28. Complete the details of Example 3.10 by showing that

$$P_{1,n} \rightarrow \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad P_{2,n} \rightarrow \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

29. Let  $A$  and  $B$  be  $m \times m$  symmetric matrices. Show that

$$\begin{aligned} \lambda_1(A+B) &\leq \lambda_1(A) + \lambda_1(B), \\ \lambda_m(A+B) &\geq \lambda_m(A) + \lambda_m(B) \end{aligned}$$

30. Prove Theorem 3.20.

31. Our proof of Theorem 3.24 utilized (3.9) of Theorem 3.17. Obtain an alternative proof of Theorem 3.24 by using (3.10) of Theorem 3.17.

32. Let  $A$  be an  $m \times m$  nonnegative definite 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)  $\lambda_{h-i+1}((F'AF)(F'BF)^{-1}) \geq \lambda_{m-i+1}(AB^{-1})$ , for  $i = 1, \dots, h$ .  
 (b)  $\min_F \lambda_1((F'AF)(F'BF)^{-1}) = \lambda_{m-h+1}(AB^{-1})$ .  
 (c)  $\min_F \lambda_h((F'AF)(F'BF)^{-1}) = \lambda_m(AB^{-1})$ .

33. Suppose  $A$  is an  $m \times m$  matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$  and associated eigenvectors  $x_1, \dots, x_m$ , while  $B$  is  $n \times n$  with eigenvalues  $\gamma_1, \dots, \gamma_n$  and eigenvectors  $y_1, \dots, y_n$ . 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

$$C = \begin{bmatrix} C_1 & (0) & \cdots & (0) \\ (0) & C_2 & \cdots & (0) \\ \vdots & \vdots & \cdots & \vdots \\ (0) & (0) & \cdots & C_r \end{bmatrix}$$

in terms of the eigenvalues and eigenvectors of the matrices  $C_1, \dots, C_r$ .

34. 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  $T'T$ .
35. Show that if  $A$  is a nonnegative definite matrix and  $a_{ii} = 0$  for some  $i$ , then  $a_{ij} = a_{ji} = 0$  for all  $j$ .
36. Suppose that  $A$  is an  $m \times m$  symmetric matrix with eigenvalues  $\lambda_1, \dots, \lambda_m$  and associated eigenvectors  $x_1, \dots, x_m$ , while  $B$  is an  $m \times m$  symmetric matrix with eigenvalues  $\gamma_1, \dots, \gamma_m$  and associated eigenvectors  $x_1, \dots, x_m$ ; 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  $AB = BA$ .
37. Suppose that  $x_1, \dots, x_r$  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,

$$P = \sum_{i=1}^r x_i x_i'$$

Let  $B$  be another  $m \times m$  symmetric matrix with its  $r$  largest eigenvalues given by  $\mu_1, \dots, \mu_r$ , where  $\mu_r > \mu_{r+1}$ , and a corresponding set of orthonormal eigenvectors given by  $y_1, \dots, y_r$ . Let  $Q$  be the total eigenprojection of  $B$  associated with the eigenvalues  $\mu_1, \dots, \mu_r$  so that

$$Q = \sum_{i=1}^r y_i y_i'$$

(a) Show that  $P = Q$  if and only if

$$\sum_{i=1}^r \{\gamma_i + \mu_i - \lambda_i(A + B)\} = 0$$

(b) Let  $X = (x_1, \dots, x_m)$ , where  $x_{r+1}, \dots, x_m$  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  $(m - r) \times (m - r)$ . Show that the converse is not true.

38. Let  $\lambda_1 \geq \dots \geq \lambda_m$  be the eigenvalues of the  $m \times m$  symmetric matrix  $A$  and let  $x_1, \dots, x_m$  be a set of corresponding orthonormal eigenvectors. For some  $k$ , define the total eigenprojection associated with the eigenvalues  $\lambda_k, \dots, \lambda_m$  as

$$P = \sum_{i=k}^m x_i x_i'$$

Show that  $\lambda_k = \dots = \lambda_m = \lambda$  if and only if

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

39. Let  $A_1, \dots, A_k$  be  $m \times m$  symmetric matrices and let  $\tau_i$  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  $\tau_i$  has multiplicity  $r$  and has  $\mathbf{x}_1, \dots, \mathbf{x}_r$  as associated eigenvectors, then

$$P \left\{ \sum_{i=1}^k (A_i - \tau_i I)^2 \right\} P = (0)$$

## CHAPTER FOUR

# Matrix Factorizations and Matrix Norms

### 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 to us 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 (1989) and Press, Flannery, Teukolsky, and Vetterling (1992).

### 2. THE SINGULAR VALUE DECOMPOSITION

The first factorization that we consider, the singular value decomposition, could be described as the most useful because this is a factorization 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  $r > 0$ , there exist orthogonal  $m \times m$  and  $n \times n$  matrices  $P$  and  $Q$ , such that  $A = PDQ'$  and  $D = P'AQ$ , where the  $m \times n$  matrix  $D$  is given by

$$\begin{array}{ll} \text{(a) } \Delta & \text{if } r = m = n, & \text{(b) } [\Delta \quad (0)] & \text{if } r = m < n, \\ \text{(c) } \begin{bmatrix} \Delta \\ (0) \end{bmatrix} & \text{if } r = n < m, & \text{(d) } \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} & \text{if } r < m, r < n, \end{array}$$

and  $\Delta$  is an  $r \times r$  diagonal matrix with positive diagonal elements. The diagonal elements of  $\Delta^2$  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  $\Delta^2$  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 3.23. Define  $\Delta$  to be the diagonal matrix whose diagonal elements are the positive square roots of the corresponding diagonal elements of  $\Delta^2$ . 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 which 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

$$\begin{aligned} P'AQ &= \begin{bmatrix} P_1'AQ_1 & P_1'AQ_2 \\ P_2'AQ_1 & P_2'AQ_2 \end{bmatrix} = \begin{bmatrix} \Delta^{-1}Q_1'A'AQ_1 & \Delta^{-1}Q_1'A'AQ_2 \\ P_2'AQ_1 & P_2'AQ_2 \end{bmatrix} \\ &= \begin{bmatrix} \Delta^{-1}\Delta^2 & \Delta^{-1}Q_1'A'(0) \\ (0) & P_2'(0) \end{bmatrix} = \begin{bmatrix} \Delta & (0) \\ (0) & (0) \end{bmatrix} \quad \square \end{aligned}$$



The diagonal elements of  $\Delta$ , 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  $A'A = QD'DQ'$ . It is important to note also that the columns of  $P$  form an orthonormal set of eigenvectors of  $AA'$  since  $AA' = PDQ'QD'P' = PDD'P'$ .

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 there exist  $m \times r$  and  $n \times r$  matrices  $P_1$  and  $Q_1$ , such that  $P_1'P_1 = Q_1'Q_1 = I_r$  and  $A = P_1\Delta Q_1'$ , where  $\Delta$  is an  $r \times r$  diagonal matrix with positive diagonal elements.

Quite a bit 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  $A$ , while 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 1.9 and its corollary, Corollary 1.9.1, which were stated as consequences of the properties of elementary transformations. It is easily verified that Theorem 1.9 and Corollary 1.9.1 also follow directly from Theorem 4.1 and Corollary 4.1.1.

**Example 4.1.** We will find a singular value decomposition for the  $4 \times 3$  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 be the columns of the  $3 \times 3$  orthogonal matrix  $Q$ . Clearly,  $\text{rank}(A) = 2$  and the two singular values of  $A$  are  $\sqrt{28}$  and  $\sqrt{12}$ . Thus, the  $4 \times 2$  matrix  $P_1$  is given by

$$P_1 = AQ_1\Delta^{-1} = \begin{bmatrix} 2 & 0 & 1 \\ 3 & -1 & 1 \\ -2 & 4 & 1 \\ 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & 1/\sqrt{3} \\ -1/\sqrt{2} & 1/\sqrt{3} \\ 0 & 1/\sqrt{3} \end{bmatrix} \begin{bmatrix} 1/\sqrt{28} & 0 \\ 0 & 1/\sqrt{12} \end{bmatrix}$$

$$= \begin{bmatrix} 1/\sqrt{14} & 1/2 \\ 2/\sqrt{14} & 1/2 \\ -3/\sqrt{14} & 1/2 \\ 0 & 1/2 \end{bmatrix}$$

The  $4 \times 2$  matrix  $P_2$  can be any matrix satisfying  $P_1'P_2 = (0)$  and  $P_2'P_2 = I_2$ ; for instance, we can take  $(1/\sqrt{12}, 1/\sqrt{12}, 1/\sqrt{12}, -3/\sqrt{12})'$  and  $(-5/\sqrt{42}, 4/\sqrt{42}, 1/\sqrt{42}, 0)'$  as the columns of  $P_2$ . Then our singular value decomposition of  $A$  is given by

$$\begin{bmatrix} 1/\sqrt{14} & 1/2 & 1/\sqrt{12} & -5/\sqrt{42} \\ 2/\sqrt{14} & 1/2 & 1/\sqrt{12} & 4/\sqrt{42} \\ -3/\sqrt{14} & 1/2 & 1/\sqrt{12} & 1/\sqrt{42} \\ 0 & 1/2 & -3/\sqrt{12} & 0 \end{bmatrix} \begin{bmatrix} \sqrt{28} & 0 & 0 \\ 0 & \sqrt{12} & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

$$\cdot \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \\ 1/\sqrt{6} & 1/\sqrt{6} & -2/\sqrt{6} \end{bmatrix},$$

or in the form of Corollary 4.1.1,

$$\begin{bmatrix} 1/\sqrt{14} & 1/2 \\ 2/\sqrt{14} & 1/2 \\ -3/\sqrt{14} & 1/2 \\ 0 & 1/2 \end{bmatrix} \begin{bmatrix} \sqrt{28} & 0 \\ 0 & \sqrt{12} \end{bmatrix} \begin{bmatrix} 1/\sqrt{2} & -1/\sqrt{2} & 0 \\ 1/\sqrt{3} & 1/\sqrt{3} & 1/\sqrt{3} \end{bmatrix}$$

Alternatively, we could have determined the matrix  $P$  by using 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 will have to 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  $x$  be an  $m \times 1$  nonnull vector. Its singular value decomposition will be of the form

$$x = Pdq,$$

where  $P$  is an  $m \times m$  orthogonal matrix,  $d$  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  $x$  is given by  $\lambda^{1/2}$ , where  $\lambda = x'x$ . If we define  $x_* = \lambda^{-1/2}x$ , note that  $x'_*x_* = 1$ , and

$$xx'x_* = xx'(\lambda^{-1/2}x) = (\lambda^{-1/2}x)x'x = \lambda x_*$$

so that  $x_*$  is a normalized eigenvector of  $xx'$  corresponding to its single positive eigenvalue  $\lambda$ . Any vector orthogonal to  $x_*$  is an eigenvector of  $xx'$  corresponding to the repeated eigenvalue 0. Thus, if we let  $d = (\lambda^{1/2}, 0, \dots, 0)'$ ,  $q = 1$ , and  $P = [x_*, p_2, \dots, p_m]$  be any orthogonal matrix with  $x_*$  as its first column, then

$$Pdq = [x_*, p_2, \dots, p_m] \begin{bmatrix} \lambda^{1/2} \\ 0 \\ \vdots \\ 0 \end{bmatrix} 1 = \lambda^{1/2}x_* = x$$

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  $2 \times 2$  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  $AA'$ , while  $A'A$  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 which illustrates an application of the singular value decomposition to least squares regression. For more discussion of this and other uses of the singular value decomposition in statistics, the reader is referred to 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, which we first discussed in Example 3.6. Suppose we have the standardized regression model

$$y = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}$$

We have seen in Example 2.15 that the least squares estimator of  $\delta_0$  is  $\bar{y}$ . The fitted model  $\hat{y} = \bar{y} \mathbf{1}_N + Z_1 \hat{\boldsymbol{\delta}}_1$  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  $y = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}$  as we did in Example 2.15 by defining  $\alpha_0 = \delta_0$ ,  $\boldsymbol{\alpha}_1 = U' \boldsymbol{\delta}_1$ , and  $W_1 = VD$ , so that  $y = \alpha_0 \mathbf{1}_N + W_1 \boldsymbol{\alpha}_1 + \boldsymbol{\epsilon}$ . 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 standard-

ized explanatory variable space, are actually confined to a  $(k - r)$ -dimensional subspace. Also, the model  $y = \alpha_0 \mathbf{1}_N + W_1 \alpha_1 + \epsilon$  simplifies to

$$y = \alpha_0 \mathbf{1}_N + W_{11} \alpha_{11} + \epsilon, \quad (4.5)$$

where  $W_{11} = V_1 D_1$  and  $\alpha_{11} = U_1' \delta_1$ , and the least squares estimator of the  $(k - r) \times 1$  vector  $\alpha_{11}$  is given by  $\hat{\alpha}_{11} = (W_{11}' W_{11})^{-1} W_{11}' y = D_1^{-1} V_1' y$ . This can be used to find a least squares estimator of  $\delta_1$  since we must have  $\hat{\alpha}_{11} = U_1' \hat{\delta}_1$ . Partitioning  $\hat{\delta}_1 = (\hat{\delta}_{11}', \hat{\delta}_{12}')'$  and  $U_1' = (U_{11}', U_{12}')$ , where  $\hat{\delta}_{11}$  is  $(k - r) \times 1$ , we obtain the relationship

$$\hat{\alpha}_{11} = U_{11}' \hat{\delta}_{11} + U_{12}' \hat{\delta}_{12}$$

Premultiplying this by  $U_{11}'^{-1}$  (if  $U_{11}$  is not nonsingular, then  $\delta_1$  and  $U_1$  can be rearranged so that it is), we find that

$$\hat{\delta}_{11} = U_{11}'^{-1} \hat{\alpha}_{11} - U_{11}'^{-1} U_{12}' \hat{\delta}_{12};$$

that is, the least squares estimator of  $\delta_1$  is not unique since  $\hat{\delta}_1 = (\hat{\delta}_{11}', \hat{\delta}_{12}')'$  is a least squares estimator for any choice of  $\hat{\delta}_{12}$ , as long as  $\hat{\delta}_{11}$  satisfies the identity given. 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  $z$ . Using a least squares estimate  $\hat{\delta}_1$  we obtain the estimate  $\hat{y} = \bar{y} + z' \hat{\delta}_1$ . This estimated response, like  $\hat{\delta}_1$ , may not be unique since, if we partition  $z$  as  $z' = (z_1', z_2')$  with  $z_1$   $(k - r) \times 1$ ,

$$\begin{aligned} \hat{y} &= \bar{y} + z' \hat{\delta}_1 = \bar{y} + z_1' \hat{\delta}_{11} + z_2' \hat{\delta}_{12} \\ &= \bar{y} + z_1' U_{11}'^{-1} \hat{\alpha}_{11} + (z_2' - z_1' U_{11}'^{-1} U_{12}') \hat{\delta}_{12} \end{aligned}$$

Thus,  $\hat{y}$  does not depend on the arbitrary  $\hat{\delta}_{12}$  and is therefore unique, only if

$$(z_2' - z_1' U_{11}'^{-1} U_{12}') = 0', \quad (4.6)$$

in which case the unique estimated value is given by  $\hat{y} = \bar{y} + z_1' U_{11}'^{-1} \hat{\alpha}_{11}$ . It is easily shown that the set of all vectors  $z = (z_1', z_2')'$  satisfying (4.6) 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 3.6, principal components regression utilizes the regression model (4.5), 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  $z \in S$  we will introduce bias into our estimate of  $y$ . The advantage of principal components regression is that this 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 4.9). 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  $S$ .

### 3. THE SPECTRAL DECOMPOSITION AND SQUARE ROOT MATRICES OF A SYMMETRIC MATRIX

The spectral decomposition of a symmetric matrix, briefly discussed in the previous chapter, is nothing more than a special case of the singular value decomposition. We summarize this result in the following theorem.

**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  $m \times m$  matrix  $A^{1/2}$  for which  $A = A^{1/2} A^{1/2}$ . If  $\Lambda$  and  $X$  are defined as in the theorem above, 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 since  $X'X = I$ ,

$$A^{1/2}A^{1/2} = X\Lambda^{1/2}X'X\Lambda^{1/2}X' = X\Lambda^{1/2}\Lambda^{1/2}X' = X\Lambda X' = A,$$

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.

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 since

$$A^{1/2}A^{1/2'} = X\Lambda^{1/2}Q'Q\Lambda^{1/2}X' = X\Lambda^{1/2}\Lambda^{1/2}X' = X\Lambda X' = A$$

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  $A$ . The following theorem establishes the existence of such a decomposition.

**Theorem 4.3.** Let  $A$  be an  $m \times m$  nonnegative definite matrix. Then there exists an  $m \times m$  lower triangular matrix  $T$  having nonnegative diagonal elements such that  $A = TT'$ . Further, if  $A$  is positive definite, 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  $m = 1$ , since 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

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

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  $t_{12}$  and a unique positive scalar  $t_{22}$  such that

$$\begin{bmatrix} A_{11} & a_{12} \\ a'_{12} & a_{22} \end{bmatrix} = \begin{bmatrix} T_{11} & \mathbf{0} \\ t'_{12} & t_{22} \end{bmatrix} \begin{bmatrix} T'_{11} & t_{12} \\ \mathbf{0}' & t_{22} \end{bmatrix} = \begin{bmatrix} T_{11}T'_{11} & T_{11}t_{12} \\ t'_{12}T'_{11} & t'_{12}t_{12} + t_{22}^2 \end{bmatrix};$$

that is, we must have  $a_{12} = T_{11}t_{12}$  and  $a_{22} = t'_{12}t_{12} + t^2_{22}$ . Since  $T_{11}$  must be nonsingular, the unique choice of  $t_{12}$  is given by  $t_{12} = T^{-1}_{11}a_{12}$ , and so  $t^2_{22}$  must satisfy

$$\begin{aligned} t^2_{22} &= a_{22} - t'_{12}t_{12} = a_{22} - a'_{12}(T^{-1}_{11})'T^{-1}_{11}a_{12} \\ &= a_{22} - a'_{12}(T_{11}T'_{11})^{-1}a_{12} = a_{22} - a'_{12}A^{-1}_{11}a_{12} \end{aligned}$$

Note that since  $A$  is positive definite,  $a_{22} - a'_{12}A^{-1}_{11}a_{12}$  will be positive since, if we let  $x = (x'_1, -1)' = (a'_{12}A^{-1}_{11}, -1)'$ , then

$$\begin{aligned} x'Ax &= x'_1A_{11}x_1 - 2x'_1a_{12} + a_{22} \\ &= a'_{12}A^{-1}_{11}A_{11}A^{-1}_{11}a_{12} - 2a'_{12}A^{-1}_{11}a_{12} + a_{22} \\ &= a_{22} - a'_{12}A^{-1}_{11}a_{12} \end{aligned}$$

Consequently, the unique  $t_{22} > 0$  is given by  $t_{22} = (a_{22} - a'_{12}A^{-1}_{11}a_{12})^{1/2}$ .  $\square$

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  $m \geq n$ . There exist an  $n \times n$  upper triangular matrix  $R$  and an  $m \times n$  matrix  $Q$  satisfying  $Q'Q = I_n$ , such that  $A = QR$ .

For a proof of Theorem 4.4 see Horn and Johnson (1985). 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  $x$  has mean vector  $\mu$  and the positive definite covariance matrix  $\Omega$ . By using a square root matrix of  $\Omega$ , we can determine a linear transformation of  $x$  so that the transformed random vector is standardized; that is, it has mean vector  $0$  and covariance matrix  $I_m$ . If we let  $\Omega^{1/2}$  be any matrix satisfying  $\Omega = \Omega^{1/2}(\Omega^{1/2})'$  and put  $z = \Omega^{-1/2}(x - \mu)$ , where  $\Omega^{-1/2} = (\Omega^{1/2})^{-1}$ , then by using (1.8) and (1.9) of Section 1.13, we find that

$$E(z) = E\{\Omega^{-1/2}(x - \mu)\} = \Omega^{-1/2}\{E(x - \mu)\} = \Omega^{-1/2}(\mu - \mu) = 0,$$

and



$$\begin{aligned}\text{var}(z) &= \text{var}\{\Omega^{-1/2}(x - \mu)\} = \Omega^{-1/2}\{\text{var}(x - \mu)\}(\Omega^{-1/2})' \\ &= \Omega^{-1/2}\{\text{var}(x)\}(\Omega^{-1/2})' = \Omega^{-1/2}\Omega(\Omega^{-1/2})' = I_m\end{aligned}$$

Since the covariance matrix of  $z$  is the identity matrix, the Euclidean distance function will give a meaningful measure of the distance between observations from this distribution. By making use of the linear transformation defined above, we can relate distances between  $z$  observations to distances between  $x$  observations. For example, the Euclidean distance between an observation  $z$  and its expected value  $\mathbf{0}$  is

$$\begin{aligned}d_1(z, \mathbf{0}) &= \{(z - \mathbf{0})'(z - \mathbf{0})\}^{1/2} = (z'z)^{1/2} \\ &= \{(x - \mu)'(\Omega^{-1/2})'\Omega^{-1/2}(x - \mu)\}^{1/2} \\ &= \{(x - \mu)'\Omega^{-1}(x - \mu)\}^{1/2} \\ &= d_\Omega(x, \mu),\end{aligned}$$

where  $d_\Omega$  is the Mahalanobis distance function defined in Section 2.2. Similarly, if  $x_1$  and  $x_2$  are two observations from the distribution of  $x$  and  $z_1$  and  $z_2$  are the corresponding transformed vectors, then  $d_1(z_1, z_2) = d_\Omega(x_1, x_2)$ . This relationship between the Mahalanobis distance and the Euclidean distance makes the construction of the Mahalanobis distance function more apparent. It 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, while the second stage simply computes the Euclidean distance for these transformed points.

**Example 4.6.** In Example 2.16, we obtained the weighted least squares estimator of  $\beta$  in the multiple regression model

$$y = X\beta + \epsilon,$$

where  $\text{var}(\epsilon) = \sigma^2 \text{diag}(c_1^2, \dots, c_N^2)$  and  $c_1^2, \dots, c_N^2$  are known constants. 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. As with weighted least squares regression, the approach here is to 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 utilizing 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

$$y_* = X_*\beta + \epsilon_*$$

where  $y_* = T^{-1}y$ ,  $X_* = T^{-1}X$ , 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  $\beta$  in the model  $y = X\beta + \epsilon$  is given by the ordinary least squares estimator of  $\beta$  in the model  $y_* = X_*\beta + \epsilon_*$  and so can be expressed as

$$\begin{aligned}\hat{\beta}_* &= (X'_*X_*)^{-1}X'_*y_* = (X'T'^{-1}T^{-1}X)^{-1}X'T'^{-1}T^{-1}y \\ &= (X'C^{-1}X)^{-1}X'C^{-1}y\end{aligned}$$

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 our next theorem, 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  $A = CC'$ , 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 there exists an  $h \times n$  matrix  $Q$  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  $h = n$  since if  $h < n$ , we can form the matrix  $B_* = [B \quad (0)]$  so that  $B_*$  is  $m \times n$  and  $B_*B_*' = BB'$ ; then proving that there exists an  $n \times n$  orthogonal matrix  $Q_*$  such that  $C = B_*Q_*$  will yield  $C = BQ$ , if we take  $Q$  to be the first  $h$  rows of  $Q_*$ .

Now since  $BB'$  is symmetric, there exists an orthogonal matrix  $X$  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  $\Lambda$  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

$$E = \begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' B = \begin{bmatrix} \Lambda^{-1/2} X_1' B \\ X_2' B \end{bmatrix} = \begin{bmatrix} E_1 \\ E_2 \end{bmatrix}, \quad (4.7)$$

$$F = \begin{bmatrix} \Lambda^{-1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix} X' C = \begin{bmatrix} \Lambda^{-1/2} X_1' C \\ X_2' C \end{bmatrix} = \begin{bmatrix} F_1 \\ F_2 \end{bmatrix}, \quad (4.8)$$

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

$$\begin{aligned} EQ &= EE_*' F_* = \begin{bmatrix} E_1 \\ (0) \end{bmatrix} [E_1' \ E_3'] \begin{bmatrix} F_1 \\ F_3 \end{bmatrix} \\ &= \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} F_1 \\ F_3 \end{bmatrix} = \begin{bmatrix} F_1 \\ (0) \end{bmatrix} = F \end{aligned}$$

But using (4.7) and (4.8),  $EQ = F$  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

$$X \begin{bmatrix} \Lambda^{1/2} & (0) \\ (0) & I_{m-r} \end{bmatrix},$$

since  $XX' = I_m$ . □

#### 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 the following definition.

**Definition 4.1.** The  $m \times m$  matrices  $A$  and  $B$  are said to be similar matrices if there exists a nonsingular matrix  $C$  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 since 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  $\Lambda$  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 diagonalizability 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 3.6.

**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  $X = (x_1, \dots, x_m)$ , where  $x_1, \dots, x_m$  are eigenvectors of  $A$  corresponding to  $\lambda_1, \dots, \lambda_m$ , then

$$X^{-1}AX = \Lambda \tag{4.9}$$

The theorem above 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  $\mu_1, \dots, \mu_h$  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  $A = X\Lambda X^{-1}$ . 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, since  $\mu_i$  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

$$(A - \mu_i I_m)x = 0$$

But any such  $x$  is an eigenvector of  $A$  corresponding to the eigenvalue  $\mu_i$ . Consequently, we can find a set of  $r_i$  linearly independent eigenvectors associated with the eigenvalue  $\mu_i$ . From Theorem 3.6, 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  $\mu_i$  for each  $i$ , will also be linearly independent. Therefore,  $A$  is diagonalizable and so the proof is complete.  $\square$

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 (4.9) immediately yields the following generalization of this result.

**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 4.7.** Let  $A$ ,  $B$ , and  $C$  be the  $2 \times 2$  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  $x_1 = (1, 2)'$  and  $x_2 = (1, -2)'$ , 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  $Bx = \lambda x = 0$  has only one linearly independent solution for  $x$ , 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  $r = 2$ , since its characteristic equation simplifies to  $(1 - \lambda)^2 = 0$ . This matrix is not diagonalizable since  $\text{rank}(C - \lambda I_2) = \text{rank}(C - I_2) = \text{rank}(B) = 1 \neq 0 = m - r$ . Any eigenvector of  $C$  is a scalar multiple of the vector  $x = (1, 0)'$ . 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.

The next result 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 and 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 2.21, we know that

$$\text{rank}(A) = m - \dim\{N(A)\},$$

where  $N(A)$  is the null space of  $A$ . But since the null space of  $A$  consists of all vectors  $x$  satisfying  $Ax = \mathbf{0}$ , 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 4.8.** We saw in Example 4.7 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 since each has only one linearly independent eigenvector associated with its single eigenvalue, which has multiplicity two. This eigenvalue is 0 for  $B$ , so

$$\text{rank}(B) = 2 - \dim\{S_B(0)\} = 2 - 1 = 1$$

On the other hand, since 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.

## 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 diagonal as is possible. We begin with the following definition.

**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 Examples 4.7 and 4.8 are both  $2 \times 2$  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 since  $J_h(\lambda)$  is a triangular matrix, its diagonal elements are its eigenvalues, and so it has the one value,  $\lambda$ , repeated  $h$  times. However, the solution to  $J_h(\lambda)x = \lambda x$  has  $x_1$  arbitrary while  $x_2 = \cdots = x_h = 0$ ; that is,  $J_h(\lambda)$  has only one linearly independent eigenvector, which is of the form  $x = (x_1, 0, \dots, 0)'$ .

We now state the Jordan decomposition theorem. For a proof of this result see Horn and Johnson (1985).

**Theorem 4.11.** Let  $A$  be an  $m \times m$  matrix. Then there exists a nonsingular matrix  $B$  such that

$$B^{-1}AB = J = \text{diag}(J_{h_1}(\lambda_1), \dots, J_{h_r}(\lambda_r)) \\ = \begin{bmatrix} J_{h_1}(\lambda_1) & (0) & \cdots & (0) \\ (0) & J_{h_2}(\lambda_2) & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & J_{h_r}(\lambda_r) \end{bmatrix},$$

where  $h_1 + \cdots + h_r = m$  and  $\lambda_1, \dots, \lambda_r$  are the not necessarily distinct eigenvalues of  $A$ .

The matrix  $J$  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  $x_i$  is an eigenvector of  $J$  corresponding to the eigenvalue  $\lambda_i$  if and only if the vector  $y_i = Bx_i$  is an eigenvector of  $A$  corresponding to  $\lambda_i$ ; for instance, if  $x_i$  satisfies  $Jx_i = \lambda_i x_i$ , then

$$Ay_i = (BJB^{-1})Bx_i = BJx_i = \lambda_i Bx_i = \lambda_i y_i$$



Thus,  $r$  also gives the number of linearly independent eigenvectors of  $A$ , and  $A$  is diagonalizable only if  $J$  is diagonal.

**Example 4.9.** Suppose that  $A$  is a  $4 \times 4$  matrix with the eigenvalue  $\lambda$  having multiplicity 4. Then  $A$  will be similar to one of the following five Jordan canonical forms:

$$\text{diag}(J_1(\lambda), J_1(\lambda), J_1(\lambda), J_1(\lambda)) = \begin{bmatrix} \lambda & 0 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix},$$

$$\text{diag}(J_2(\lambda), J_1(\lambda), J_1(\lambda)) = \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix},$$

$$\text{diag}(J_3(\lambda), J_1(\lambda)) = \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 1 & 0 \\ 0 & 0 & \lambda & 0 \\ 0 & 0 & 0 & \lambda \end{bmatrix},$$

$$\text{diag}(J_2(\lambda), J_2(\lambda)) = \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 0 & 0 \\ 0 & 0 & \lambda & 1 \\ 0 & 0 & 0 & \lambda \end{bmatrix},$$

$$J_4(\lambda) = \begin{bmatrix} \lambda & 1 & 0 & 0 \\ 0 & \lambda & 1 & 0 \\ 0 & 0 & \lambda & 1 \\ 0 & 0 & 0 & \lambda \end{bmatrix}$$

The first form given is diagonal so this corresponds to the case in which  $A$  has four linearly independent eigenvectors associated with the eigenvalue  $\lambda$ . 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.

## 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 the following theorem.

**Theorem 4.12.** Let  $A$  be an  $m \times m$  matrix. Then there exists an  $m \times m$  unitary matrix  $X$  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  $y_1$  be an eigenvector of  $A$  corresponding to  $\lambda_1$  and normalized so that  $y_1^*y_1 = 1$ . Let  $Y$  be any  $m \times m$  unitary matrix having  $y_1$  as its first column. Writing  $Y$  in partitioned form as  $Y = [y_1 \ Y_2]$ , we see that, since  $Ay_1 = \lambda_1 y_1$  and  $Y_2^*y_1 = 0$ ,

$$\begin{aligned} Y^*AY &= \begin{bmatrix} y_1^*Ay_1 & y_1^*AY_2 \\ Y_2^*Ay_1 & Y_2^*AY_2 \end{bmatrix} = \begin{bmatrix} \lambda_1 y_1^*y_1 & y_1^*AY_2 \\ \lambda_1 Y_2^*y_1 & Y_2^*AY_2 \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & y_1^*AY_2 \\ 0 & B \end{bmatrix}, \end{aligned}$$

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, since 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  $\lambda_2$  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 there exists a unitary matrix  $W$  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 since  $W$  is. If we let  $X = YU$ , then  $X$  is also unitary and

$$\begin{aligned} X^*AX &= U^*Y^*AYU = \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W^* \end{bmatrix} \begin{bmatrix} \lambda_1 & y_1^*AY_2 \\ \mathbf{0} & B \end{bmatrix} \begin{bmatrix} 1 & \mathbf{0}' \\ \mathbf{0} & W \end{bmatrix} \\ &= \begin{bmatrix} \lambda_1 & y_1^*AY_2W \\ \mathbf{0} & W^*BW \end{bmatrix} = \begin{bmatrix} \lambda_1 & y_1^*AY_2W \\ \mathbf{0} & T_2 \end{bmatrix}, \end{aligned}$$

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 there exist corresponding real eigenvectors. 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 there exists an  $m \times m$  orthogonal matrix  $X$  such that  $X'AX = T$ , where  $T$  is an upper triangular matrix.

**Example 4.10.** Consider the  $3 \times 3$  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  $2 \times 2$  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  $A$  of Example 4.10.

In Chapter 3, by utilizing 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, while 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

$$|A| = |XTX^*| = |X^*X| |T| = |T| = \prod_{i=1}^m \lambda_i,$$

since  $|X^*X| = 1$  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

$$\operatorname{tr}(A) = \operatorname{tr}(XTX^*) = \operatorname{tr}(X^*XT) = \operatorname{tr}(T) = \sum_{i=1}^m \lambda_i$$

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  $\operatorname{rank}(A) \geq r$ .

*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 since the determinant of a triangular matrix is the product of its diagonal elements, so we must have  $\operatorname{rank}(T) \geq r$ . The result then follows from the fact that since  $X$  is unitary, it must be nonsingular, so

$$\operatorname{rank}(A) = \operatorname{rank}(XTX^*) = \operatorname{rank}(T) \geq r \quad \square$$

## 7. THE SIMULTANEOUS DIAGONALIZATION OF TWO SYMMETRIC MATRICES

We have already discussed in Section 3.7 one manner in which two symmetric matrices can be simultaneously diagonalized. We restate this result in the following theorem.

**Theorem 4.14.** Let  $A$  and  $B$  be  $m \times m$  symmetric matrices, with  $A$  being nonnegative definite and  $B$ , 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 there exists a nonsingular matrix  $C$  such that

$$CAC' = \Lambda, \quad CBC' = I_m$$

**Example 4.11.** One application of the simultaneous diagonalization described in the theorem above is in a multivariate analysis commonly referred to as canonical variate analysis (see Krzanowski, 1988 or Mardia, Kent, and Bibby, 1979). This analysis involves data from the multivariate one-way classification model discussed in Example 3.14, 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  $y_{i1}, \dots, y_{in_i}$ . The model is

$$y_{ij} = \mu_i + \epsilon_{ij},$$

where  $\mu_i$  is an  $m \times 1$  vector of constants and  $\epsilon_{ij} \sim N_m(\mathbf{0}, \Omega)$ . In Example 3.14, we saw how the matrices

$$B = \sum_{i=1}^k n_i (\bar{y}_i - \bar{y})(\bar{y}_i - \bar{y})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)(y_{ij} - \bar{y}_i)',$$

where

$$\bar{y}_i = \sum_{j=1}^{n_i} \frac{y_{ij}}{n_i}, \quad \bar{y} = \sum_{i=1}^k \frac{n_i \bar{y}_i}{n}, \quad n = \sum_{i=1}^k n_i,$$

could be used to test the hypothesis,  $H_0: \mu_1 = \dots = \mu_k$ . 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  $\mu_1, \dots, \mu_k$  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  $B$ ,

$$\Phi = \sum_{i=1}^k n_i (\mu_i - \mu)(\mu_i - \mu)',$$

where  $\mu = \Sigma n_i \mu_i / n$ , will have rank  $r$ ; in fact, the eigenvectors of  $\Phi$  corresponding to its positive eigenvalues will span this  $r$ -dimensional space. Thus, a plot of the projections of  $\mu_1, \dots, \mu_k$  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 since Euclidean distance would not be appropriate. This difficulty can be resolved by analyzing the transformed data  $\Omega^{-1/2} y_{ij}$ , where  $\Omega^{-1/2} \Omega^{-1/2} = \Omega^{-1}$ , since  $\Omega^{-1/2} y_{ij} \sim N_m(\Omega^{-1/2} \mu_i, I_m)$ . Thus, we would plot the projections of  $\Omega^{-1/2} \mu_1, \dots, \Omega^{-1/2} \mu_k$  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  $\Lambda_1$  is an  $r \times r$  diagonal matrix, then we could simply plot the vectors  $P_1' \Omega^{-1/2} \mu_1, \dots, P_1' \Omega^{-1/2} \mu_k$  in  $R^r$ . The  $r$  components of the vector  $v_i = P_1' \Omega^{-1/2} \mu_i$  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  $\Phi$  and  $\Omega$ , since if  $C' = (C_1', C_2')$  satisfies

$$\begin{bmatrix} C_1 \\ C_2 \end{bmatrix} \Phi [C_1' \quad C_2'] = \begin{bmatrix} \Lambda_1 & (0) \\ (0) & (0) \end{bmatrix}, \quad \begin{bmatrix} C_1 \\ C_2 \end{bmatrix} \Omega [C_1' \quad C_2'] = \begin{bmatrix} I_r & (0) \\ (0) & I_{m-r} \end{bmatrix},$$

then we can take  $C_1 = P_1' \Omega^{-1/2}$ . When  $\mu_1, \dots, \mu_k$  are unknown, the canonical variate means can be estimated by the sample canonical variate means, which are computed using the samples means  $\bar{y}_1, \dots, \bar{y}_k$  and the corresponding simultaneous diagonalization of  $B$  and  $W$ .

The matrix  $C$  that diagonalizes  $A$  and  $B$  in Theorem 4.14 is nonsingular 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 9, is what we consider next; we would like to know whether or not there exists an orthogonal matrix that diagonalizes both  $A$  and  $B$ . The following result gives a necessary and sufficient condition for such an orthogonal matrix to exist.

**Theorem 4.15.** Suppose that  $A$  and  $B$  are  $m \times m$  symmetric matrices. Then there exists an orthogonal matrix  $P$  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  $P'AP = \Lambda_1$  and  $P'BP = \Lambda_2$ , where  $\Lambda_1$  and

$\Lambda_2$  are diagonal matrices. Then since  $\Lambda_1$  and  $\Lambda_2$  are diagonal matrices, clearly  $\Lambda_1\Lambda_2 = \Lambda_2\Lambda_1$ , so we have

$$AB = P\Lambda_1P'P\Lambda_2P' = P\Lambda_1\Lambda_2P' = P\Lambda_2\Lambda_1P' = P\Lambda_2P'P\Lambda_1P' = BA$$

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  $\mu_1, \dots, \mu_h$  be the distinct values of the eigenvalues of  $A$  having multiplicities  $r_1, \dots, r_h$ , respectively. Since  $A$  is symmetric there exists an orthogonal matrix  $Q$  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

$$C = Q'BQ = \begin{bmatrix} C_{11} & C_{12} & \cdots & C_{1h} \\ C_{21} & C_{22} & \cdots & C_{2h} \\ \vdots & \vdots & \ddots & \vdots \\ C_{h1} & C_{h2} & \cdots & C_{hh} \end{bmatrix},$$

where  $C_{ij}$  is  $r_i \times r_j$ . Note that since  $AB = BA$ , we must have

$$\Lambda_1 C = Q'AQQ'BQ = Q'ABQ = Q'BAQ = Q'BQQ'AQ = C\Lambda_1$$

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 since  $C$  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  $\Delta_i$  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. □

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 since  $A$  and  $B$  are symmetric,  $(AB)' = B'A' = BA$ , and so  $AB = BA$  if and only if  $AB$  is symmetric. The previous theorem easily generalizes to a collection of symmetric matrices.

**Theorem 4.16.** Let  $A_1, \dots, A_k$  be  $m \times m$  symmetric matrices. Then there exists an orthogonal matrix  $P$  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 4.16 is a special case of the following result. The proof, which is similar to that given for Theorem 4.15, is left as an exercise.

**Theorem 4.17.** Suppose that each of the  $m \times m$  matrices  $A_1, \dots, A_k$  is diagonalizable. Then there exists a nonsingular matrix  $X$  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  $(i, j)$ .

## 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) since 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.

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 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

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

while 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 utilizes the eigenvalues of  $A^*A$ ; in particular, if  $\mu_1, \dots, \mu_m$  are the eigenvalues of  $A^*A$ , then the spectral norm is given by

$$\|A\|_2 = \max_{1 \leq i \leq m} \sqrt{\mu_i}$$

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 4.18.** 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. The following result shows that  $\rho(A)$  actually serves as a lower bound for any matrix norm of  $A$ .

**Theorem 4.19.** For any  $m \times m$  matrix  $A$  and any matrix norm  $\|A\|$ ,  $\rho(A) \leq \|A\|$ .

*Proof.* Suppose that  $\lambda$  is an eigenvalue of  $A$  for which  $|\lambda| = \rho(A)$ , and let  $x$  be a corresponding eigenvector, so that  $Ax = \lambda x$ . Then  $x\mathbf{1}'_m$  is an  $m \times m$  matrix satisfying  $Ax\mathbf{1}'_m = \lambda x\mathbf{1}'_m$ , and so using properties (c) and (e) of matrix norms, we find that

$$\rho(A)\|x\mathbf{1}'_m\| = |\lambda|\|x\mathbf{1}'_m\| = \|\lambda x\mathbf{1}'_m\| = \|Ax\mathbf{1}'_m\| \leq \|A\|\|x\mathbf{1}'_m\|.$$

The result now follows by dividing the equation above by  $\|x\mathbf{1}'_m\|$ . □

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 4.20.** 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  $c > 0$ , define 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

$$\sum_{j=1}^{i-1} |c^{-(i-j)} t_{ji}| < \epsilon,$$

for each  $i$ . In this case, since  $|\lambda_i| \leq \rho(A)$ , 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

$$\|B\|_{A,\epsilon} = \|(XD_c^{-1})^{-1} B (XD_c^{-1})\|_1$$

The result now follows from Theorem 4.18 and the fact that

$$\|A\|_{A,\epsilon} = \|(XD_c^{-1})^{-1} A (XD_c^{-1})\|_1 = \|D_c T D_c^{-1}\|_1 \quad \square$$

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  $x_1, x_2, \dots$  converges to  $x$ , with respect to the vector norm  $\|x\|$ , if  $\|x_k - x\| \rightarrow 0$  as  $k \rightarrow \infty$ . The following result indicates that the actual choice of a norm is not important. For a proof of this result, see Horn and Johnson (1985).

**Theorem 4.21.** Let  $\|x\|_a$  and  $\|x\|_b$  be any two vector norms defined on any  $m \times 1$  vector  $x$ . If  $x_1, x_2, \dots$  is a sequence of  $m \times 1$  vectors, then  $x_k$  converges to  $x$ , as  $k \rightarrow \infty$ , with respect to  $\|x\|_a$  if and only if  $x_k$  converges to  $x$ , as  $k \rightarrow \infty$ , with respect to  $\|x\|_b$ .

Since the first four conditions of a matrix norm are the conditions of a vector norm, the previous theorem immediately leads to the following.

**Corollary 4.21.1.** 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 4.22.** 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  $k \rightarrow \infty$ , since  $\|A\| < 1$ . Thus,  $A^k$  converges to  $(0)$  with respect to the norm  $\|A\|$ . But by Corollary 4.21.1,  $A^k$  also converges to  $(0)$  with respect to the matrix norm (see Problem 4.37)

$$\|A\|_* = m \left( \max_{1 \leq i, j \leq m} |a_{ij}| \right)$$

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 4.23.** 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  $\lambda$ , we must also have  $\lambda^k \mathbf{x} \rightarrow \mathbf{0}$ , since  $A^k \mathbf{x} = \lambda^k \mathbf{x}$ . This can only happen if  $|\lambda| < 1$ , and so  $\rho(A) < 1$ , since  $\lambda$  was an arbitrary eigenvalue of  $A$ . On the other hand, if  $\rho(A) < 1$ , then we know from Theorem 4.20 that there is a matrix norm satisfying  $\|A\| < 1$ . Hence, it follows from Theorem 4.22 that  $A^k \rightarrow (0)$ .  $\square$

Our final result shows that the spectral radius of  $A$  is the limit of a particular sequence that can be computed from any matrix norm.

**Theorem 4.24.** 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.*  $\lambda$  is an eigenvalue of  $A$  if and only if  $\lambda^k$  is an eigenvalue of  $A^k$ . Further,  $|\lambda|^k = |\lambda^k|$ , so  $\rho(A)^k = \rho(A^k)$ . This, along with Theorem 4.19, 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  $\epsilon > 0$ , there exists an integer  $N_\epsilon$  such that  $\|A^k\|^{1/k} < \rho(A) + \epsilon$  for all  $k > N_\epsilon$ . But this is the same as showing that there exists an integer  $N_\epsilon$  such that for all  $k > N_\epsilon$ ,  $\|A^k\| < \{\rho(A) + \epsilon\}^k$ , or equivalently,

$$\|B^k\| < 1, \tag{4.10}$$

where  $B = \{\rho(A) + \epsilon\}^{-1}A$ . Now (4.10) follows immediately from Theorem 4.23 since

$$\rho(B) = \frac{\rho(A)}{\rho(A) + \epsilon} < 1 \quad \square$$

## PROBLEMS

1. Obtain a singular value decomposition for the matrix

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

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  $\alpha A$  compare to those of  $A$ ?
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. 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 nonzero singular values of  $AB$  are not the same as those of  $BA$ .
5. Let  $A$  be an  $m \times n$  matrix having rank  $r$  and singular values  $\mu_1, \dots, \mu_r$ . 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.

6. Find a singular value decomposition for the vector  $x = (1, 5, 7, 5)'$ .
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  $y$ .
8. 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  $BAB$ .
9. If  $t$  is an estimator of  $\theta$ , then the mean squared error (MSE) of  $t$  is defined by

$$\text{MSE}(t) = \text{var}(t) + \{E(t) - \theta\}^2$$

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 which has the standardized explanatory variables at the values given in the  $k \times 1$  vector  $z$ . Let  $\hat{y} = \bar{y} + z'(Z_1'Z_1)^{-1}Z_1'y$  be the estimate obtained using ordinary least squares, while  $\tilde{y} = \bar{y} + z'U_1D_1^{-1}V_1'y$  is the estimate obtained using principal components regression. Assume throughout that  $\epsilon \sim N_N(\mathbf{0}, \sigma^2 I_N)$ .

(a) Show that if the vector  $\mathbf{v} = (v_1, \dots, v_N)'$  satisfies  $\mathbf{z}' = \mathbf{v}'D\mathbf{U}'$ , 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  $\text{MSE}(\tilde{y}) < \text{MSE}(\hat{y})$ ?

10. 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  $y = \delta_0 \mathbf{1}_N + Z_1 \boldsymbol{\delta}_1 + \boldsymbol{\epsilon}$ , and obtain the fitted values  $\hat{y} = \hat{\delta}_0 \mathbf{1}_N + Z_1 \hat{\boldsymbol{\delta}}_1$ .
- (b) Compute the singular value decomposition of  $Z_1$ . 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  $x_2 = 4$ .

11. Consider the  $3 \times 3$  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 a symmetric square root matrix for  $A$ .  
 (c) Find a nonsymmetric square root matrix for  $A$ .

12. Use the spectral decomposition theorem to prove Theorem 4.5.

13. Find a  $3 \times 2$  matrix  $T$  such that  $TT' = A$ , where

$$A = \begin{bmatrix} 5 & 4 & 0 \\ 4 & 5 & 3 \\ 0 & 3 & 5 \end{bmatrix}$$

14. Suppose  $x \sim N_3(\mathbf{0}, \Omega)$ , where

$$\Omega = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 2 & 1 \\ 1 & 1 & 2 \end{bmatrix}$$

Find a  $3 \times 3$  matrix  $A$  such that the components of  $z = Ax$  are independently distributed.

15. Let the matrices  $A$ ,  $B$ , and  $C$  be given by

$$A = \begin{bmatrix} 1 & 2 & 5 \\ 2 & 1 & 4 \\ -1 & 1 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 1 & -1 \\ -2 & 2 & 2 \\ -1 & 3 & 1 \end{bmatrix}, \quad C = \begin{bmatrix} 2 & 1 & -1 \\ 2 & 5 & 3 \\ -2 & -1 & 1 \end{bmatrix}$$

- (a) Which of these matrices are diagonalizable?  
 (b) Which of these matrices have their rank equal to the number of nonzero eigenvalues?

16. 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.

17. Find a  $4 \times 4$  matrix  $A$  having eigenvalues 0 and 1 with multiplicities 3 and 1, respectively, such that
- the rank of  $A$  is 1,
  - the rank of  $A$  is 2,
  - the rank of  $A$  is 3.
18. Repeat Example 4.9 for  $5 \times 5$  matrices; that is, obtain a collection of  $5 \times 5$  matrices in Jordan canonical form such that every  $5 \times 5$  matrix having the eigenvalue  $\lambda$  with multiplicity 5 is similar to one of the matrices in this set.
19. Consider the  $6 \times 6$  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.

- Find the eigenvalues of  $J$  and their multiplicities.
  - Find the eigenspaces of  $J$ .
20. An  $m \times m$  matrix  $B$  is said to be nilpotent if  $B^k = (0)$  for some positive integer  $k$ .
- Show that  $J_h(\lambda) = \lambda I_h + B_h$ , where  $B_h$  is nilpotent. In particular, show that  $B_h^h = (0)$ .
  - 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  $B^h = (0)$ ?
  - 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.
21. Let  $A$  be an  $m \times m$  nilpotent matrix. In Problem 3.27, 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  $A^h = (0)$ .
22. 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  $\text{rank}(A^2) = \text{rank}(A)$ .

23. Suppose that  $\lambda$  is an eigenvalue of  $A$  with multiplicity  $r$ . Show that there are  $r$  linearly independent eigenvectors of  $A$  corresponding to  $\lambda$  if and only if  $\text{rank}(A - \lambda I) = \text{rank}\{(A - \lambda I)^2\}$ .
24. Let  $A$  and  $B$  be  $m \times m$  matrices. Suppose that there exists an  $m \times m$  unitary matrix  $X$  such  $X^*AX$  and  $X^*BX$  are both upper triangular matrices. Show then that the eigenvalues of  $AB - BA$  are all equal to 0.
25. 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  $1 \leq i \leq r+1$ ,  $1 \leq j \leq r+1$ .
26. Use the Schur decomposition of a matrix  $A$  and the result of the previous exercise to prove the Cayley–Hamilton theorem given as Theorem 3.7; that is, if  $\lambda_1, \dots, \lambda_m$  are the eigenvalues of  $A$ , show that

$$(A - \lambda_1 I)(A - \lambda_2 I) \cdots (A - \lambda_m I) = (0).$$

27. Obtain a Schur decomposition for the matrix  $C$  given in Problem 15.
28. Repeat Problem 27 by obtaining a different Schur decomposition of  $C$ .
29. Let  $A$  be  $m \times n$ , with  $m \leq n$ . Show that there exist an  $m \times m$  nonnegative definite matrix  $B$  and an  $m \times n$  matrix  $H$  such that  $HH' = I_m$  and  $A = BH$ .
30. 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 there exists a nonsingular matrix  $X$  such that both  $X^{-1}AX$  and  $X^{-1}BX$  are diagonal matrices. This proves Theorem 4.17 when  $k = 2$ .

31. 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 30?
32. 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  $\mu_1, \dots, \mu_m$ .

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  $(1, \dots, m)$ .

33. The following is a generalization of Theorem 4.14 to arbitrary nonnegative definite matrices. Let  $A$  and  $B$  be  $m \times m$  nonnegative definite matrices with  $\text{rank}(A) = r \leq s = \text{rank}(B)$ . Show that there exists a nonsingular matrix  $C$  such that

$$CAC' = \begin{bmatrix} D_1 & (0) \\ (0) & (0) \end{bmatrix}, \quad CBC' = \begin{bmatrix} D_2 & (0) \\ (0) & (0) \end{bmatrix},$$

where  $D_1$  and  $D_2$  are  $t \times t$  diagonal matrices,  $t \geq s$ , and

$$D_1 = \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix}$$

34. 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.
35. 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.
36. Show that the functions,  $\|A\|_E$ ,  $\|A\|_1$ ,  $\|A\|_\infty$ , and  $\|A\|_2$  given in Section 4.8 are, in fact, matrix norms.
37. 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.

38. Prove Theorem 4.18.
39. 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.
40. Show that if for some matrix norm  $\|I_m - A\| < 1$  then  $A$  is a nonsingular matrix.

## PROBLEMS

41. Consider the  $2 \times 2$  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  $\|A\| < 1$ .

42. Let  $A$  be an  $m \times m$  matrix. Show that if for some matrix norm  $\|A\| < 1$  then the matrix  $I_m - A$  has an inverse and

$$(I_m - A)^{-1} = I_m + \sum_{k=1}^{\infty} A^k$$

43. 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. Use this 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  $2 \times 2$  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  $A\mathbf{x} = \mathbf{c}$ .

44. Suppose that  $A$  is an  $m \times m$  matrix. Show that there exist an  $m \times m$  lower triangular matrix  $L$ , an  $m \times m$  upper triangular matrix  $U$ , and  $m \times m$  permutation matrices  $P$  and  $Q$ , such that  $A = PLUQ$ .

45. 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  $A = LDL'$ .

## CHAPTER FIVE

# Generalized Inverses

### 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 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

$$Ax = c,$$

where  $A$  is an  $m \times n$  matrix of constants,  $c$  is an  $m \times 1$  vector of constants, and  $x$  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  $x = A^{-1}c$ ; that is, the system has a solution, the solution is unique, and it is given by  $x = A^{-1}c$ . 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 the next chapter 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  $x$  which has a mean vector of zero and covariance matrix  $\Omega$ . A useful transformation in some situations is one that transforms  $x$  to another random vector,  $z$ , having the identity matrix as its covariance matrix. For instance, in Chapter 9 we will see that if  $z$  has a normal distribution, then the sum of squares of the components of  $z$ , that is  $z'z$ , 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  $z = T'x$  will have  $I_m$  as its covariance matrix. Then

$$z'z = x'(T')'T'x = x'(TT')x = x'\Omega^{-1}x$$

This, of course, will be possible only if  $\Omega$  is positive definite. If  $\Omega$  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  $z$  is defined by  $z = B'x$ ,

$$\text{var}(z) = \begin{bmatrix} I_r & (0) \\ (0) & (0) \end{bmatrix},$$

and  $z'z = x'Ax$ . We will see later that  $A$  is a generalized inverse of  $\Omega$  and  $z'z$  still has a chi-squared distribution if  $z$  has a normal distribution.

## 2. THE MOORE-PENROSE GENERALIZED INVERSE

A very 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 which 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 the following theorem.

**Theorem 5.1.** Corresponding to each  $m \times n$  matrix  $A$ , there exists one and only one  $n \times m$  matrix  $A^+$  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 there exist  $m \times r$  and  $n \times r$  matrices  $P$  and  $Q$  such that  $P'P = Q'Q = I_r$  and

$$A = P\Delta Q',$$

where  $\Delta$  is a diagonal matrix with positive diagonal elements. Define  $A^+ = Q\Delta^{-1}P'$ , and note that

$$AA^+A = P\Delta Q'Q\Delta^{-1}P'P\Delta Q' = P\Delta\Delta^{-1}\Delta Q' = P\Delta Q' = A$$

$$A^+AA^+ = Q\Delta^{-1}P'P\Delta Q'Q\Delta^{-1}P' = Q\Delta^{-1}\Delta\Delta^{-1}P' = Q\Delta^{-1}P' = A^+$$

$$AA^+ = P\Delta Q'Q\Delta^{-1}P' = PP' \quad \text{is symmetric}$$

$$A^+A = Q\Delta^{-1}P'P\Delta Q' = QQ' \quad \text{is symmetric}$$

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

$$AB = (AB)' = B'A' = B'(ACA)' = B'A'(AC)' = (AB)'AC = ABAC = AC,$$

and

$$BA = (BA)' = A'B' = (ACA)'B' = (CA)'A'B' = CA(BA)' = CABA = CA$$

Now using these two identities, 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 very 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, is the original definition given by Moore (1935). This definition utilizes the concept of projection matrices that were 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  $x \in R^m$ ,  $P_Sx$  gives the orthogonal projection of  $x$  onto  $S$ , while  $x - P_Sx$  is the component of  $x$  orthogonal to  $S$ ; further, this unique matrix



$P_S$  is given by  $x_1 x_1' + \cdots + x_r x_r'$ , where  $\{x_1, \dots, x_r\}$  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

$$(a) \quad AA^+ = P_{R(A)}, \quad (b) \quad A^+A = P_{R(A^+)},$$

where  $P_{R(A)}$  and  $P_{R(A^+)}$  are the projection matrices of the range spaces of  $A$  and  $A^+$ , respectively.

The equivalence of Definitions 5.1 and 5.2 is not immediately obvious. Consequently, we will establish it in the next theorem.

**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 since by definition, a projection matrix is symmetric, while (5.1) and (5.2) follow since 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 2.19 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).  $\square$

### 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, while in some of the subsequent sections, we will look at some more specialized results. First, we have the following theorem.

**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

$$A'A(A'A)^+ = A'AA^+A^{+'} = A'(A^+(AA^+))' = A'(A^+AA^+)' = A'A^{+'} = (A^+A)',$$

and  $A^+A$  must be symmetric by definition, so condition (5.3) is satisfied for  $(A'A)^+ = A^+A^{+'}$ . Likewise condition (5.4) holds since

$$(A'A)^+A'A = A^+A^{+'}A'A = A^+(AA^+)'A = A^+AA^+A = A^+A$$

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) since  $\text{rank}(A) = 2$ . 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

$$A^+ = A'(AA')^{-1} = \frac{1}{14} \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix} \begin{bmatrix} 5 & -4 \\ -4 & 6 \end{bmatrix} = \frac{1}{14} \begin{bmatrix} -3 & 8 \\ 6 & -2 \\ 5 & -4 \end{bmatrix}$$

Our next result 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  $P_{R(A')} = A^+A$ .

Our next result 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 since

$$AAA = A^2A = AA = A^2 = A,$$

while conditions (5.3) and (5.4) hold because

$$(AA)' = A'A' = AA \quad \square$$

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  $\Lambda$  be the  $m \times m$  diagonal matrix  $\text{diag}(\lambda_1, \dots, \lambda_m)$ . Then the Moore-Penrose inverse  $\Lambda^+$  of  $\Lambda$ , 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  $\lambda_i$ 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  $\Lambda$  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

$$A = [X_1 \ X_2] \begin{bmatrix} \Lambda_1 & (0) \\ (0) & (0) \end{bmatrix} \begin{bmatrix} X_1' \\ X_2' \end{bmatrix} = X_1 \Lambda_1 X_1',$$

and similarly the expression above for  $A^+$  reduces to  $A^+ = X_1 \Lambda_1^{-1} X_1'$ . Thus, since  $X_1' X_1 = I_r$ , we have

$$AA^+ = X_1 \Lambda_1 X_1' X_1 \Lambda_1^{-1} X_1' = X_1 \Lambda_1 \Lambda_1^{-1} X_1' = X_1 X_1',$$

which is clearly symmetric, so condition (5.3) is satisfied. Similarly,  $A^+A = X_1 X_1'$  and so (5.4) also holds. Conditions (5.1) and (5.2) hold since

$$AA^+A = (AA^+)A = X_1 X_1' X_1 \Lambda_1 X_1' = X_1 \Lambda_1 X_1' = A$$

and

$$A^+AA^+ = A^+(AA^+) = X_1 \Lambda_1^{-1} X_1' X_1 X_1' = X_1 \Lambda_1^{-1} X_1' = A^+,$$

and so the proof is complete. □

**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 have

$$P_{R(X)} = XX^+ = X(X'X)^+X' \quad (5.7)$$

**Example 5.3.** We will utilize (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 exercise. Using this we 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 since  $P_{R(X)}$  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 the normalized eigenvectors  $z_1 = (1/\sqrt{2}, -1/\sqrt{2}, 0)'$  and  $z_2 = (1/\sqrt{6}, 1/\sqrt{6}, 2/\sqrt{6})'$  corresponding to its two positive eigenvalues. Thus, if we let  $Z = (z_1, z_2)$ , 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 9; here we will look at a simple illustration. A common situation encountered in inferential statistics is one in which one has a sample statistic,  $t \sim N_m(\theta, \Omega)$ , and it is desired to determine whether or not the  $m \times 1$  parameter vector  $\theta = \mathbf{0}$ ; formally, we want to test the null hypothesis  $H_0: \theta = \mathbf{0}$  versus the alternative hypothesis  $H_1: \theta \neq \mathbf{0}$ . One approach to this problem, if  $\Omega$  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

$$\text{var}(\mathbf{u}) = T^{-1} \{\text{var}(\mathbf{t})\} T'^{-1} = T^{-1}(TT')T'^{-1} = \mathbf{I}_m,$$

so  $\mathbf{u} \sim N_m(T^{-1}\boldsymbol{\theta}, \mathbf{I}_m)$ . 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  $\Omega$  is positive semidefinite, the construction of  $v_1$  above can be generalized by using the Moore-Penrose inverse of  $\Omega$ . 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  $\Lambda_1$  are defined as in the proof of Theorem 5.7, then  $\mathbf{w} = \Lambda_1^{-1/2}X_1'\mathbf{t} \sim N_r(\Lambda_1^{-1/2}X_1'\boldsymbol{\theta}, \mathbf{I}_r)$ , since

$$\text{var}(\mathbf{w}) = \Lambda_1^{-1/2}X_1' \{\text{var}(\mathbf{t})\} X_1\Lambda_1^{-1/2} = \Lambda_1^{-1/2}X_1'(X_1\Lambda_1X_1')X_1\Lambda_1^{-1/2} = \mathbf{I}_r$$

Thus, since 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.

#### 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  $2 \times 1$  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

$$(\mathbf{a}'\mathbf{b})^+ = (1)^+ = 1 \neq \mathbf{b}^+\mathbf{a}^+ = 0.5$$

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

$$(AA^+)^+ = AA^+ = (A^+)^+A^+$$

The next theorem gives yet another situation.

**Theorem 5.8.** Let  $A$  be an  $m \times n$  matrix, while  $P$  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 the theorem above.

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  $AB$ .  $\square$

While 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  $C$ ,

$$C^+CC' = (C^+C)'C' = C'C^+C' = C'C'^+C' = C' \quad (5.9)$$

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, due to condition (5.4), is equivalent to

$$B^+A^+AB = (AB)^+(AB) = P_{R((AB)^+)} \tag{5.10}$$

The final equality in (5.10) follows from the definition of the Moore-Penrose inverse in terms 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 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

$$ABB^+A^+ = (AB)(AB)^+ = P_{R(AB)} \tag{5.11}$$

But 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  $(AB)^+ = B^+A^+$ . Using this 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

$$ABB'(I - A^+A)BB'A' = (0)$$

Note that since  $D = (I - A^+A)$  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,

$$(I - A^+A)BB'A' = (0),$$

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  $B' A' A A'$ , 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  $(AB)^+$  which 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  $AB = A_1 B_1$  and  $(A_1 B_1)^+ = B_1^+ A_1^+$ . The result, due to Cline (1964a), is given in the next theorem.

**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  $A_1 = AB_1 B_1^+$ , then  $AB = A_1 B_1$  and  $(AB)^+ = B_1^+ A_1^+$ .

*Proof.* Note that

$$AB = AA^+ AB = AB_1 = AB_1 B_1^+ B_1 = A_1 B_1,$$

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

$$A^+ A_1 = A^+ AB_1 B_1^+ = A^+ A(A^+ AB)B_1^+ = A^+ ABB_1^+ = B_1 B_1^+, \quad (5.12)$$

and

$$A_1^+ A_1 = A_1^+ AB_1 B_1^+ = A_1^+ A(B_1 B_1^+ B_1)B_1^+ = A_1^+ A_1 B_1 B_1^+ \quad (5.13)$$

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

$$A_1' = (AB_1 B_1^+)' = (AB_1 B_1^+ B_1 B_1^+)' = (A_1 B_1 B_1^+)' = B_1 B_1^+ A_1',$$

and then postmultiplying this identity by  $A_1 B_1$ .  $\square$

Note that in Theorem 5.11,  $B$  was transformed to  $B_1$  by the projection matrix of the range space of  $A^+$ , while  $A$  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  $AB = A_1B_1$ . 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  $B$ , a  $p \times n$  matrix. If we define  $B_1 = A^+AB$  and  $A_1 = ABB^+$ , then  $(AB)^+ = B_1^+A_1^+$ .

### 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  $C = (I_m - UU^+)V$ ,  $M = \{I_{n_2} + (I_{n_2} - C^+C)V'U^+U^+V(I_{n_2} - C^+C)\}^{-1}$ , and  $W = (I_{n_2} - C^+C)MV'U^+U^+(I_m - VC^+)$ .

The proof of Theorem 5.13, which is rather lengthy, will be omitted. The interested reader should refer to Cline (1964b), Boullion and Odell (1971), or Pringle and Rayner (1971). The proofs of the following consequences of Theorem 5.13 can also be found in these references.

**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

$$(a) \ A^+ = \begin{bmatrix} U^+ - U^+VKV'U^+U^+ \\ C^+ + KV'U^+U^+ \end{bmatrix}$$

if and only if  $C^+CV'U^+U^+V = (0)$ ,

$$(b) \ A^+ = \begin{bmatrix} U^+ - U^+VKV'U^+U^+ \\ KV'U^+U^+ \end{bmatrix}$$

if and only if  $C = (0)$ ,

$$(c) \ A^+ = \begin{bmatrix} U^+ - U^+VC^+ \\ C^+ \end{bmatrix}$$

if and only if  $C^+CV'U^+U^+V = V'U^+U^+V$ ,

$$(d) A^+ = \begin{bmatrix} U^+ \\ V^+ \end{bmatrix}$$

if and only if  $U^+V = (0)$ .

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

$$A^+ = \begin{bmatrix} A_{11}^+ & (0) & \cdots & (0) \\ (0) & A_{22}^+ & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & A_{rr}^+ \end{bmatrix}$$

## 6. THE MOORE–PENROSE INVERSE OF A SUM

Theorem 1.7 gave an expression for  $(A+CBD)^{-1}$ , when the matrices  $A$  and  $B$  are both 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 utilize 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

$$(UU' + VV')^+ = (I_m - C^+V')U^+K U^+(I_m - VC^+) + (CC')^+,$$

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 the following theorem.

**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

$$VU^+ = VU'(UU')^+ = \{(UU')^+UV'\}' = (0)$$

Similarly, we have  $V^+U = (0)$  and  $UV^+ = (0)$ . As a result,

$$(U + V)(U^+ + V^+) = UU^+ + VV^+ \tag{5.14}$$

$$(U^+ + V^+)(U + V) = U^+U + V^+V, \tag{5.15}$$

which are both symmetric, so that conditions (5.3) and (5.4) are satisfied. Post-multiplying equation (5.14) by  $(U + V)$  and (5.15) 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

$$(U_1 + \dots + U_k)^+ = U_1^+ + \dots + U_k^+$$

## 7. THE CONTINUITY OF THE MOORE-PENROSE INVERSE

It is very useful to establish the continuity of a function since 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  $A$ . But before doing this, 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, due to the continuity of sums and the continuity of scalar products, we immediately have the following.

**Theorem 5.18.** 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

$$A^{-1} = |A|^{-1} A_{\#}, \quad (5.16)$$

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, due to the continuity of the determinant function,  $|A_i| \rightarrow |A|$ , and so there must exist an  $N$  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 of  $A_i$ , then  $A_{i\#} \rightarrow A_{\#}$  as  $i \rightarrow \infty$ . As a result, equation (5.16) has allowed us to establish the following.

**Theorem 5.19.** Let  $A$  be an  $m \times m$  nonsingular matrix. Then the inverse of  $A$ ,  $A^{-1}$ , is a continuous function of the elements of  $A$ .

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  $2 \times 2$  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  $\text{rank}(A) = 1$ , while  $\text{rank}(A_i) = 2$  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 since its (1, 1)th element,  $i$ , converges to  $\infty$ . 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 the example above; that is, as  $A_i$  gets closer to  $A$ ,  $A_i^+$  will get closer to  $A^+$ . This continuity property of  $A^+$  is summarized below. A proof of this important result can be found in Penrose (1955) or Campbell and Meyer (1979).

**Theorem 5.20.** Let  $A$  be an  $m \times n$  matrix and  $A_1, A_2, \dots$  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 there exists an integer  $N$  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  $\theta$  if  $t$  converges in probability to  $\theta$ ; 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  $\theta$ , 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  $\Omega$  be an  $m \times m$  positive semidefinite covariance matrix having rank  $r < m$ . Suppose that the elements of the matrix  $\Omega$  are unknown and are, therefore, to be estimated. Suppose, in addition, that our sample estimate of  $\Omega$ , 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  $\Omega$ ; that is, each element of  $\hat{\Omega}$  is a consistent estimator of the corresponding element of  $\Omega$ . However, since  $\text{rank}(\Omega) = r < m$ ,  $\hat{\Omega}^+$  is not a consistent estimator of  $\Omega^+$ . 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  $\Lambda$  are converging to zero, while the  $m - r$  largest diagonal elements of  $\Lambda^{-1}$  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  $\Omega$  having rank  $r$ . For example, if  $\hat{\Omega}$  has eigenvalues  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  and corresponding normalized eigenvectors  $x_1, \dots, x_m$  and  $P_r$  is the eigenprojection

$$P_r = \sum_{i=1}^r x_i x_i'$$

then

$$\hat{\Omega}_* = P_r \hat{\Omega} P_r = \sum_{i=1}^r \lambda_i x_i x_i'$$

will be an estimator of  $\Omega$  having rank of  $r$ . It can be shown then that, due to the continuity of eigenprojections,  $\hat{\Omega}_*$  is also a consistent estimator of  $\Omega$ . More importantly, since  $\text{rank}(\hat{\Omega}_*) = \text{rank}(\Omega) = r$ , Theorem 5.20 guarantees that  $\hat{\Omega}_*^+$  is a consistent estimator of  $\Omega^+$ .

## 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 utilizing some of the four conditions (5.1)–(5.4) or, for sim-

plicity, 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 the next chapter. Discussion of other useful  $\{i_1, \dots, i_r\}$ -inverses can be found in Ben-Israel and Greville (1974), Campbell and Meyer (1979), and Rao and Mitra (1971).

In the next chapter, 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  $A$  makes use of 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 5.21.** 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  $\Delta$  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  $B$  above 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, since  $PP' = I_m$ ,  $QQ' = I_n$ , any generalized inverse  $B$ , of  $A$ , can be expressed as

$$\begin{aligned} B &= QQ'BP' = Q \begin{bmatrix} Q'_1 \\ Q'_2 \end{bmatrix} B[P_1 \ P_2]P' \\ &= Q \begin{bmatrix} Q'_1BP_1 & Q'_1BP_2 \\ Q'_2BP_1 & Q'_2BP_2 \end{bmatrix} P', \end{aligned}$$

which is in the required form if we can show that  $Q'_1BP_1 = \Delta^{-1}$ . Since  $B$  is a generalized inverse of  $A$ ,  $ABA = A$ , or equivalently,  $P'ABAQ = P'AQ$ . Writing this last identity in partitioned form and equating the (1, 1)th submatrices on both sides, we find that

$$\Delta Q'_1BP_1\Delta = \Delta,$$

from which it immediately follows that  $Q'_1BP_1 = \Delta^{-1}$ , and so the proof is complete.  $\square$

**Example 5.8.** The  $4 \times 3$  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 5.21, 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

$$G = \begin{bmatrix} 1/\sqrt{6} & 0 \end{bmatrix},$$

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, while the Moore–Penrose inverse has its rank equal to that of  $A$ , which is 2.

The following theorem summarizes some of the basic properties of  $\{1\}$ -inverses.

**Theorem 5.22.** 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  $\alpha$  is a nonzero scalar,  $\alpha^{-1}A^-$  is a generalized inverse of  $\alpha A$ ,
- (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  $A^-A = I_n$ .

*Proof.* Properties (a)–(d) are easily proven by simply verifying that the one condition of a generalized inverse holds. To prove (e), note that since  $A = AA^{-}A$ , we can use Theorem 2.10 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). □

**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 of the generalized inverses of a matrix  $A$  can be expressed in terms of any one particular generalized inverse. This relationship is given below.

**Theorem 5.23.** 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  $AA^{-1}A = A$ ,

$$\begin{aligned} A(A^{-1} + C - A^{-1}ACAA^{-1})A &= AA^{-1}A + ACA - AA^{-1}ACAA^{-1}A \\ &= A + ACA - ACA = A, \end{aligned}$$

so  $A^{-1} + C - A^{-1}ACAA^{-1}$  is a generalized inverse of  $A$  regardless of the choice of  $A^{-1}$  and  $C$ . Now let  $B$  be any generalized inverse of  $A$  and define  $C = B - A^{-1}$ . Then, since  $ABA = A$ , we have

$$\begin{aligned} A^{-1} + C - A^{-1}ACAA^{-1} &= A^{-1} + (B - A^{-1}) - A^{-1}A(B - A^{-1})AA^{-1} \\ &= B - A^{-1}ABAA^{-1} + A^{-1}AA^{-1}AA^{-1} \\ &= B - A^{-1}AA^{-1} + A^{-1}AA^{-1} = B, \end{aligned}$$

and so the proof is complete. □

We will find the following result useful in a later chapter.

**Theorem 5.24.** 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  $C(ABC)^{-1}A$  is a generalized inverse of  $B$ .

*Proof.* Our proof follows that of Srivastava and Khatri (1979). Using Theorem 2.10, 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

$$\text{rank}(AB) = \text{rank}(BC) = \text{rank}(B) = \text{rank}(ABC) \quad (5.17)$$

Using Theorem 2.12 along with the identity

$$A(BC)\{I_q - (ABC)^{-1}ABC\} = (0),$$

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 (5.17). But this can be true only if

$$BC\{I_q - (ABC)^-ABC\} = \{I_q - BC(ABC)^-A\}B(C) = (0)$$

Again applying Theorem 2.12, 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 (5.17). This implies that

$$\{I_q - BC(ABC)^-A\}B = B - B\{C(ABC)^-A\}B = (0),$$

and so the result follows. □

We will see in the next chapter 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 the  $\{1, 3\}$ -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 5.22 also apply to  $A^L$ . Some additional properties of least squares inverses are given below.

**Theorem 5.25.** 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. To prove (b) first note that

$$\begin{aligned} A(A'A)^- A'A &= AA^+ A(A'A)^- A'A = (AA^+)' A(A'A)^- A'A = A^{+'} A' A(A'A)^- A'A \\ &= A^{+'} A'A = (AA^+)' A = AA^+ A = A, \end{aligned}$$

so that  $(A'A)^- A'$  satisfies condition 1. To verify that condition 3 holds, observe that

$$\begin{aligned} A(A'A)^- A' &= A(A'A)^- A' A^{+'} A' = A(A'A)^- A' (AA^+)' \\ &= A(A'A)^- A' AA^+ = AA^+, \end{aligned}$$

where the last equality uses the identity,  $A(A'A)^- A'A = A$ , just proven. Thus, the symmetry of  $A(A'A)^- A'$  follows from the symmetry of  $AA^+$ .  $\square$

## 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  $[B \ c]$ , 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  $A_j = (a_1, \dots, a_j)$ , so that  $A_j$  is the  $m \times j$  matrix containing the first  $j$  columns of  $A$ . Greville has shown that if we write  $A_j = [A_{j-1} \ a_j]$ , then

$$A_j^+ = \begin{bmatrix} A_{j-1}^+ & -d_j b_j' \\ & b_j' \end{bmatrix}, \tag{5.18}$$

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  $A_2 = [\mathbf{a}_1 \ \mathbf{a}_2] = [A_1 \ \mathbf{a}_2]$ . We find that

$$A_1^+ = (\mathbf{a}'_1 \mathbf{a}_1)^{-1} \mathbf{a}'_1 = \left( \frac{1}{3}, \frac{1}{3}, \frac{1}{3} \right),$$

$$\mathbf{d}_2 = A_1^+ \mathbf{a}_2 = \frac{1}{3},$$

$$\mathbf{c}_2 = \mathbf{a}_2 - A_1 \mathbf{d}_2 = \mathbf{a}_2 - \frac{1}{3} \mathbf{a}_1 = \frac{1}{3} \begin{bmatrix} 2 \\ -4 \\ 2 \end{bmatrix}$$

Since  $\mathbf{c}_2 \neq \mathbf{0}$ , we get

$$\mathbf{b}'_2 = \mathbf{c}_2^+ = (\mathbf{c}'_2 \mathbf{c}_2)^{-1} \mathbf{c}'_2 = \frac{1}{4} [1, \ -2, \ 1],$$

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  $A_3 = [A_2 \ \mathbf{a}_3]$  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

$$\mathbf{b}'_3 = (1 + \mathbf{d}'_3 \mathbf{d}_3)^{-1} \mathbf{d}'_3 A_2^+ = (1 + 2)^{-1} [1 \quad 1] \frac{1}{4} \begin{bmatrix} 1 & 2 & 1 \\ 1 & -2 & 1 \end{bmatrix}$$

$$= \frac{1}{6} [1 \quad 0 \quad 1],$$

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 5.26.** Let  $A$  be an  $m \times m$  matrix. Then there exists a nonsingular  $m \times m$  matrix  $C$  such that  $CA = H$ , where  $H$  is in Hermite form.

The proof of the next result will be left to the reader as an exercise.

**Theorem 5.27.** 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 5.26 will be a generalized inverse of  $A$ .

**Theorem 5.28.** 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 5.27 we know that  $H$  is idempotent and so

$$CACA = H^2 = H = CA$$

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  $3 \times 3$  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, while the remaining elements in the first column are all equal to zero. This can be achieved via the matrix equation  $C_1A = A_1$  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, while 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 5.28 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

$$C = C_2 C_1 = \frac{1}{6} \begin{bmatrix} 2 & 0 & 1 \\ 2 & -1 & 0 \\ 12 & -12 & 12 \end{bmatrix}$$

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. The following result indicates how such an extension is possible.

**Theorem 5.29.** 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 5.28 that  $C$  is a generalized inverse of  $A_*$ . Hence,  $A_* C A_* = A_*$ . Simplifying the left-hand side of this identity, we find that

$$\begin{aligned} A_* C A_* &= \begin{bmatrix} A \\ (0) \end{bmatrix} [C_1 \ C_2] \begin{bmatrix} A \\ (0) \end{bmatrix} \\ &= \begin{bmatrix} AC_1 & AC_2 \\ (0) & (0) \end{bmatrix} \begin{bmatrix} A \\ (0) \end{bmatrix} = \begin{bmatrix} AC_1 A \\ (0) \end{bmatrix} \end{aligned}$$

Equating this to  $A_*$ , we get  $AC_1 A = A$ , 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

$$A_* = [A \quad 0] = \begin{bmatrix} 1 & 1 & 2 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 1 & 2 & 0 \\ 2 & 0 & 2 & 0 \end{bmatrix}$$

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 5.25(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

$$A^L = (A'A)^{-1}A' = \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix}$$

### PROBLEMS

1. Prove results (a)–(d) of Theorem 5.3.
2. Use Theorem 5.3(h) to find the Moore–Penrose inverse of

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

3. Find the Moore–Penrose inverse of the vector

$$a = \begin{bmatrix} 2 \\ 1 \\ 3 \\ 2 \end{bmatrix}$$

4. Provide the proofs for (f)–(j) of Theorem 5.3.
5. Prove Theorem 5.6.
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.

7. Consider the matrix

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



## PROBLEMS

- (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  $A$ .
8. Let  $A$  be an  $m \times n$  matrix with  $\text{rank}(A) = 1$ . Show that  $A^+ = c^{-1}A'$ , where  $c = \text{tr}(A'A)$ .
9. Let  $x$  and  $y$  be  $m \times 1$  vectors, and let  $\mathbf{1}$  be the  $m \times 1$  vector with each element equal to one. Obtain expressions for the Moore–Penrose inverses of
- (a)  $\mathbf{1}\mathbf{1}'$       (b)  $I_m - m^{-1}\mathbf{1}\mathbf{1}'$       (c)  $xx'$       (d)  $xy'$ .
10. 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.
11. 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)  $A(A'A)^+A'A = AA'(AA')^+A = A$ .
12. 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.
13. Let  $A$  be an  $m \times m$  symmetric matrix having rank  $r$ . Show that if  $A$  has one nonzero eigenvalue  $\lambda$  of multiplicity  $r$ , then  $A^+ = \lambda^{-2}A$ .
14. 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
- $$AB(AB)^+ = AA^+$$
15. 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.
16. 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  $A^+A + B^+B = I_m$ .
17. 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 eigenvec-

tors 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  $(AB)^+ = B^+A^+$ .

18. Prove Theorem 5.8.

19. Prove (b)–(d) of Theorem 5.10.

20. For each case below use Theorem 5.10 to determine whether  $(AB)^+ = B^+A^+$ .

$$(a) \quad A = \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 2 \end{bmatrix}$$

$$(b) \quad A = \begin{bmatrix} 1 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}, \quad B = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix}$$

21. Let  $A$  be an  $m \times n$  matrix and  $B$  be an  $n \times m$  matrix. Show that  $(AB)^+ = B^+A^+$  if  $A'ABB' = BB'A'A$ .

22. Prove Theorem 5.14.

23. Find the Moore–Penrose inverse of the matrix

$$A = \begin{bmatrix} 2 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 0 & 1 & 2 & 0 \\ 0 & 0 & 0 & 0 & 4 \end{bmatrix}$$

24. Use Corollary 5.13.1(d) to find the Moore–Penrose inverse of the matrix  $A = [U \ V]$ , where

$$U = \begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}, \quad V = \begin{bmatrix} 1 & -2 \\ -1 & 1 \\ 0 & 1 \end{bmatrix}$$

25. Use Corollary 5.13.1(c) to find the Moore–Penrose inverse of the matrix  $A = [U \ V]$ , where

$$U = \begin{bmatrix} 1 & 1 \\ 1 & -1 \\ 1 & 0 \\ 1 & 0 \\ 0 & 0 \end{bmatrix}, \quad V = \begin{bmatrix} 2 & 2 \\ 2 & 0 \\ -1 & 0 \\ 1 & -2 \\ 0 & 1 \end{bmatrix}$$

26. 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  $A = wx' + yz'$ .

27. Find a generalized inverse, different from the Moore–Penrose inverse, of the vector given in Problem 3.
28. Consider the diagonal matrix  $A = \text{diag}(0, 2, 3)$ .
- Find a generalized inverse of  $A$  having rank of 2.
  - Find a generalized inverse of  $A$  that has rank of 3 and is diagonal.
  - Find a generalized inverse of  $A$  that is not diagonal.
29. 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  $\text{rank}(B) = n$ .
30. 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.
31. 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  $PAQ$ .
32. Let  $(A'A)^-$  be any generalized inverse of the matrix  $A'A$ , where  $A$  is  $m \times n$ . Establish the following.
- $A(A'A)^-A'$  does not depend on the choice of  $(A'A)^-$ .
  - $A(A'A)^-A'$  is symmetric even if  $(A'A)^-$  is not symmetric.
33. 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  $A(A'A)^-A' = A_1(A_1'A_1)^-A_1'$ .

34. Use the recursive procedure described in Section 9 to obtain the Moore–Penrose inverse of the matrix.

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

35. Find a generalized inverse of the matrix  $A$  in the previous exercise by finding a nonsingular matrix that transforms it into a matrix having Hermite form.
36. Find a generalized inverse of the matrix

$$A = \begin{bmatrix} 1 & -1 & -2 & 1 \\ -2 & 4 & 3 & -2 \\ 1 & 1 & -3 & 1 \end{bmatrix}$$

37. Find a least squares inverse for the matrix  $A$  given in the previous exercise.
38. It was shown in Theorem 5.28 that a generalized inverse of  $A$  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  $A$ .
39. Prove Theorem 5.27.
40. Penrose (1956) obtained the following recursive method for calculating the Moore–Penrose 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.

41. Let  $\lambda$  be the largest eigenvalue of  $AA'$ , where  $A$  is an  $m \times n$  matrix. Let  $\alpha$  be any constant satisfying  $0 < \alpha < 2/\lambda$  and define  $X_1 = \alpha A'$ . Ben-Israel

## PROBLEMS

(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  $\lambda$  does not need to be computed since we must have

$$\frac{2}{\text{tr}(AA')} < \frac{2}{\lambda}$$

42. Use the results of Section 5 to obtain the expression given in (5.18) for the Moore–Penrose inverse of the matrix  $A_j = [A_{j-1} \quad a_j]$ .

## CHAPTER SIX

# Systems of Linear Equations

### 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

$$Ax = c, \quad (6.1)$$

where  $A$  is an  $m \times n$  matrix of constants,  $c$  is an  $m \times 1$  vector of constants, and  $x$  is an  $n \times 1$  vector of variables for which solutions are needed. 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 conclude the chapter by taking a look at the special application of finding least squares solutions to (6.1), when an exact solution does not exist.

### 2. CONSISTENCY OF A SYSTEM OF EQUATIONS

In this section, we will obtain necessary and sufficient conditions for the existence of a vector  $x$  satisfying equation (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  $c$  is in the column space of  $A$  or, equivalently, that the rank of the augmented matrix  $[A \ c]$  is the same as the rank of  $A$ .

**Theorem 6.1.** The system of equations,  $Ax = c$ , is consistent if and only if  $\text{rank}([A \ c]) = \text{rank}(A)$ .

*Proof.* If  $a_1, \dots, a_n$  are the columns of  $A$ , then the equation  $Ax = c$  can be written as

$$Ax = [a_1 \cdots a_n] \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = \sum_{i=1}^n x_i a_i = c$$

Clearly, this holds for some  $x$  if and only if  $c$  is a linear combination of the columns of  $A$ , in which case  $\text{rank}[A \ c] = \text{rank}(A)$ .  $\square$

**Example 6.1.** Consider the system of equations which has

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix}, \quad c = \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix}$$

Clearly, the rank of  $A$  is 2 while

$$|[A \ c]| = \begin{vmatrix} 1 & 2 & 1 \\ 2 & 1 & 5 \\ 1 & 0 & 3 \end{vmatrix} = 0,$$

so that the rank of  $[A \ c]$  is also 2. Thus, we know from Theorem 6.1 that the system of equations  $Ax = c$  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. Our next result gives an alternative necessary and sufficient condition for consistency utilizing a generalized inverse,  $A^-$ , of  $A$ . An obvious consequence of this result is that when the system  $Ax = c$  is consistent, then a solution will be given by  $x = A^-c$ .

**Theorem 6.2.** The system of equations  $Ax = c$  is consistent if and only if for some generalized inverse,  $A^-$ , of  $A$ ,  $AA^-c = c$ .

*Proof.* First, suppose that the system is consistent and  $x_*$  is a solution, so that  $c = Ax_*$ . Premultiplying this identity by  $AA^-$ , where  $A^-$  is any generalized inverse of  $A$ , yields

$$AA^-c = AA^-Ax_* = Ax_* = c,$$

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, since  $\mathbf{x}_* = A^{-}\mathbf{c}$  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,  $AA_1\mathbf{c} = \mathbf{c}$ . Then  $A_2$  satisfies the same condition since

$$AA_2\mathbf{c} = AA_2(AA_1\mathbf{c}) = (AA_2A)A_1\mathbf{c} = AA_1\mathbf{c} = \mathbf{c}.$$

Thus, in applying Theorem 6.2, one will need to check the given condition for only one generalized inverse of  $A$ , and it doesn't matter which generalized inverse is used. In particular, we can use the Moore–Penrose inverse  $A^+$ , of  $A$ .

The following results involve some special cases regarding the matrix  $A$ .

**Corollary 6.2.1.** If  $A$  is an  $m \times m$  nonsingular matrix and  $\mathbf{c}$  is an  $m \times 1$  vector of constants, then the system  $A\mathbf{x} = \mathbf{c}$  is consistent.

**Corollary 6.2.2.** If the  $m \times n$  matrix  $A$  has rank equal to  $m$ , then the system  $A\mathbf{x} = \mathbf{c}$  is consistent.

*Proof.* Since  $A$  has full row rank, it follows from Theorem 5.22(f) that  $AA^{-} = I_m$ . As a result,  $AA^{-}\mathbf{c} = \mathbf{c}$ , and so from Theorem 6.2, the system must be consistent.  $\square$

**Example 6.2.** Consider the system of equations  $A\mathbf{x} = \mathbf{c}$ , where

$$A = \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix}, \quad \mathbf{c} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix}$$

A generalized inverse of the transpose of  $A$  was given in Example 5.12. Using this, we find that

$$\begin{aligned} AA^{-}\mathbf{c} &= \begin{bmatrix} 1 & 1 & 1 & 2 \\ 1 & 0 & 1 & 0 \\ 2 & 1 & 2 & 2 \end{bmatrix} \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \\ &= \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 1 & 1 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} \end{aligned}$$

Since this is  $\mathbf{c}$ , the system of equations is consistent, and a solution is given by



$$A^{-}c = \begin{bmatrix} 0 & 1 & -1 \\ 1 & -1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 3 \\ 2 \\ 5 \end{bmatrix} = \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix}$$

The system of linear equations  $Ax = c$  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 the following theorem.

**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 equation (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 since

$$AX_*B = AA^{-}CB^{-}B = C \quad \square$$

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 upon the choices of generalized inverses for  $A$  and  $B$ .

### 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  $\mathbf{c} = \mathbf{0}$ , our system of equations may have many solutions. The following theorem gives a general expression for all solutions to the system.

**Theorem 6.4.** Suppose that  $A\mathbf{x} = \mathbf{c}$  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  $\mathbf{y}$ ,

$$\mathbf{x}_y = A^- \mathbf{c} + (\mathbf{I}_n - A^- A) \mathbf{y} \quad (6.3)$$

is a solution, and for any solution,  $\mathbf{x}_*$ , there exists a vector  $\mathbf{y}$  such that  $\mathbf{x}_* = \mathbf{x}_y$ .

*Proof.* Since  $A\mathbf{x} = \mathbf{c}$  is a consistent system of equations, we know from Theorem 6.2 that  $AA^- \mathbf{c} = \mathbf{c}$ , and so

$$\begin{aligned} A\mathbf{x}_y &= AA^- \mathbf{c} + A(\mathbf{I}_n - A^- A) \mathbf{y} \\ &= \mathbf{c} + (A - AA^- A) \mathbf{y} = \mathbf{c}, \end{aligned}$$

since  $AA^- A = A$ . 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,

$$A^- \mathbf{c} + (\mathbf{I}_n - A^- A) \mathbf{x}_* = A^- \mathbf{c} + \mathbf{x}_* - A^- A\mathbf{x}_* = \mathbf{x}_*,$$

so that  $\mathbf{x}_* = \mathbf{x}_{\mathbf{x}_*}$ . 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  $\mathbf{x}_*$ , there exists a generalized inverse  $B$  such that  $\mathbf{x}_* = B\mathbf{c}$ .

*Proof.* Theorem 6.4 was not dependent upon 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 since  $\mathbf{c} \neq \mathbf{0}$ , 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

$$\begin{aligned} \mathbf{x}_y &= A^- \mathbf{c} + (\mathbf{I}_n - A^- A) \mathbf{y} = A^- \mathbf{c} + (\mathbf{I}_n - A^- A) C \mathbf{c} \\ &= A^- \mathbf{c} + C \mathbf{c} - A^- A C \mathbf{c} = A^- \mathbf{c} + C \mathbf{c} - A^- A C A A^- \mathbf{c} \\ &= (A^- + C - A^- A C A A^-) \mathbf{c} \end{aligned}$$

But it follows from Theorem 5.23 that  $A^- + C - A^- A C A A^-$  is a generalized inverse of  $A$  for any choice of the  $n \times m$  matrix  $C$  and so the proof is complete.  $\square$

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,  $X_*$ , there exists a matrix  $Y$  such that  $X_* = X_Y$ .

**Example 6.3.** For the consistent system of equations discussed in Example 6.2, we have

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

Consequently, a general solution to this system of equations is given by

$$\begin{aligned} \mathbf{x}_y &= A^- \mathbf{c} + (\mathbf{I}_n - A^- A) \mathbf{y} \\ &= \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 & 1 & 1 & 2 \\ 0 & 0 & 0 & -2 \\ -2 & -1 & -1 & -2 \\ 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} \\ &= \begin{bmatrix} -3 + 2y_1 + y_2 + y_3 + 2y_4 \\ 1 - 2y_4 \\ 5 - 2y_1 - y_2 - y_3 - 2y_4 \\ y_4 \end{bmatrix} \end{aligned}$$

where  $\mathbf{y}$  is an arbitrary  $4 \times 1$  vector.

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  $y$ ?

**Theorem 6.6.** If  $Ax = c$  is a consistent system of equations, then the solution  $x_* = A^-c$  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  $x_* = A^-c$  is a unique solution if and only if  $x_y = x_*$  for all choices of  $y$ , where  $x_y$  is as defined in (6.3). In other words, the solution is unique if and only if

$$(I_n - A^-A)y = 0$$

for all  $y$ , and clearly this is equivalent to the condition  $(I_n - A^-A) = (0)$  or  $A^-A = I_n$ .  $\square$

We saw in Theorem 5.22(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.

**Corollary 6.6.1.** Suppose that  $Ax = c$  is a consistent system of equations. Then the solution  $x_* = A^-c$  is a unique solution if and only if  $\text{rank}(A) = n$ .

**Example 6.4.** We saw in Example 6.1 that the system of equations  $Ax = c$ , where

$$A = \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix}, \quad c = \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix},$$

is consistent. The Moore–Penrose inverse of the transpose of  $A$  was obtained in Example 5.1. Using this, we find that

$$A^+A = \frac{1}{14} \begin{bmatrix} -3 & 6 & 5 \\ 8 & -2 & -4 \end{bmatrix} \begin{bmatrix} 1 & 2 \\ 2 & 1 \\ 1 & 0 \end{bmatrix} = \frac{1}{14} \begin{bmatrix} 14 & 0 \\ 0 & 14 \end{bmatrix} = I_2$$

Thus, the system of equations  $Ax = c$  has the unique solution given by

$$A^+c = \frac{1}{14} \begin{bmatrix} -3 & 6 & 5 \\ 8 & -2 & -4 \end{bmatrix} \begin{bmatrix} 1 \\ 5 \\ 3 \end{bmatrix} = \frac{1}{14} \begin{bmatrix} 42 \\ -14 \end{bmatrix} = \begin{bmatrix} 3 \\ -1 \end{bmatrix}$$

Suppose that a system of linear equations has more than one solution, and let  $x_1$  and  $x_2$  be two different solutions. Then, since  $Ax_i = c$  for  $i = 1$  and  $2$ , it follows that for any scalar  $\alpha$

$$A\{\alpha x_1 + (1 - \alpha)x_2\} = \alpha Ax_1 + (1 - \alpha)Ax_2 = \alpha c + (1 - \alpha)c = c$$

Thus,  $x = \{\alpha x_1 + (1 - \alpha)x_2\}$  is also a solution. Since  $\alpha$  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  $c \neq 0$  must be between 1 and  $n$ ; that is, there exists a set of linearly independent solutions,  $\{x_1, \dots, x_r\}$ , such that every solution can be expressed as a linear combination of the solutions,  $x_1, \dots, x_r$ . In other words, any solution  $x$  can be written as  $x = \alpha_1 x_1 + \dots + \alpha_r x_r$ , for some coefficients,  $\alpha_1, \dots, \alpha_r$ . Note that since  $Ax_i = c$  for each  $i$ , we must have

$$Ax = A\left(\sum_{i=1}^r \alpha_i x_i\right) = \sum_{i=1}^r \alpha_i Ax_i = \sum_{i=1}^r \alpha_i c = \left(\sum_{i=1}^r \alpha_i\right)c,$$

and so if  $x$  is a solution, the coefficients must satisfy the identity,  $\alpha_1 + \dots + \alpha_r = 1$ . Our next result tells us exactly how to determine this number of linearly independent solutions  $r$  when  $c \neq 0$ . We will delay the discussion of the situation in which  $c = 0$  until the next section.

**Theorem 6.7.** Suppose that the system  $Ax = c$  is consistent, where  $A$  is  $m \times n$  and  $c \neq 0$ . 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,  $x_0 = A^+c$ ,  $x_{e_1} = A^+c + (I_n - A^+A)e_1, \dots, x_{e_n} = A^+c + (I_n - A^+A)e_n$ , where, as usual,  $e_i$  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 since for any  $y = (y_1, \dots, y_n)'$ ,

$$x_y = A^+c + (I_n - A^+A)y = \left(1 - \sum_{i=1}^n y_i\right)x_0 + \sum_{i=1}^n y_i x_{e_i}$$

Thus, if we define the  $n \times (n+1)$  matrix  $X = (x_0, x_{e_1}, \dots, x_{e_n})$ , 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  $B = (A^+c, I_n - A^+A)$  and

$$C = \begin{bmatrix} 1 & \mathbf{1}'_n \\ \mathbf{0} & I_n \end{bmatrix}$$

Clearly,  $C$  is nonsingular since it is lower triangular and the product of its diagonal elements is 1. Consequently, from Theorem 1.8, we know that  $\text{rank}(X) = \text{rank}(B)$ . Note also that

$$\begin{aligned} (I_n - A^+A)'A^+c &= (I_n - A^+A)A^+c = (A^+ - A^+AA^+)c \\ &= (A^+ - A^+)c = \mathbf{0}, \end{aligned}$$

so that the first column of  $B$  is orthogonal to the remaining columns. This implies that

$$\text{rank}(B) = \text{rank}(A^+c) + \text{rank}(I_n - A^+A) = 1 + \text{rank}(I_n - A^+A),$$

since the consistency condition,  $AA^+c = c$  and  $c \neq \mathbf{0}$  guarantee that  $A^+c \neq \mathbf{0}$ . All that remains is to show that  $\text{rank}(I_n - A^+A) = n - \text{rank}(A)$ . Now since  $A^+A$  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  $x_0 = A^+c$  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  $x_0$ , with the remaining solutions given by  $x_y$  for  $r - 1$  different choices of  $y \neq \mathbf{0}$ . This statement is not dependent upon the choice of  $A^+$  as the generalized inverse in (6.3), since  $A^-c$  and  $(I_n - A^-A)y$  are linearly independent regardless of the choice of  $A^-$  if  $c \neq \mathbf{0}, y \neq \mathbf{0}$ . The proof of this linear independence is left as an exercise.

**Example 6.5.** We saw that the system of equations  $Ax = c$  of Examples 6.2 and 6.3 has the set of solutions consisting of all vectors of the form

$$x_y = A^-c + (I_4 - A^-A)y = \begin{bmatrix} -3 + 2y_1 + y_2 + y_3 + 2y_4 \\ 1 - 2y_4 \\ 5 - 2y_1 - y_2 - y_3 - 2y_4 \\ y_4 \end{bmatrix}$$

Since the last row of the  $3 \times 4$  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  $y$  vector. For instance, since  $A^{-}c$  and  $(I_4 - A^{-}A)y$  are linearly independent, the three solutions

$$A^{-}c, A^{-}c + (I_4 - A^{-}A)_{\cdot i}, A^{-}c + (I_4 - A^{-}A)_{\cdot j}$$

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  $Ax = c$  are given by

$$A^{-}c = \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix}, \quad A^{-}c + (I_4 - A^{-}A)_{\cdot 1} = \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ 0 \\ -2 \\ 0 \end{bmatrix} = \begin{bmatrix} -1 \\ 1 \\ 3 \\ 0 \end{bmatrix},$$

$$A^{-}c + (I_4 - A^{-}A)_{\cdot 4} = \begin{bmatrix} -3 \\ 1 \\ 5 \\ 0 \end{bmatrix} + \begin{bmatrix} 2 \\ -2 \\ -2 \\ 1 \end{bmatrix} = \begin{bmatrix} -1 \\ -1 \\ 3 \\ 1 \end{bmatrix}$$

#### 4. HOMOGENEOUS SYSTEMS OF EQUATIONS

The system of equations  $Ax = c$  is called a nonhomogeneous system of equations when  $c \neq 0$ , while  $Ax = 0$  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 since it will always have the trivial solution,  $x = 0$ . 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  $A \neq I_n$ , 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  $\alpha_1, \dots, \alpha_r$ , since

$$Ax = A \left( \sum_{i=1}^r \alpha_i x_i \right) = \sum_{i=1}^r \alpha_i Ax_i = \sum_{i=1}^r \alpha_i \mathbf{0} = \mathbf{0}$$

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  $y$ ,

$$x_y = (I_n - A^{-1}A)y$$

is a solution to the system  $Ax = \mathbf{0}$ , and for any solution,  $x_*$ , there exists a vector  $y$  such that  $x_* = x_y$ . The following result 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  $x_*$  be any solution to the system of equations  $Ax = c$ . Then

- (a) if  $x_{\#}$  is a solution to the system  $Ax_{\#} = \mathbf{0}$ ,  $x = x_* + x_{\#}$  is a solution of  $Ax = c$ , and
- (b) for any solution  $x$  to the equation  $Ax = c$ , there exists a solution  $x_{\#}$  to the equation  $Ax = \mathbf{0}$  such that  $x = x_* + x_{\#}$ .

*Proof.* Note that if  $x_{\#}$  is as defined in (a), then

$$A(x_* + x_{\#}) = Ax_* + Ax_{\#} = c + \mathbf{0} = c,$$



and so  $x = x_* + x_\#$  is a solution to  $Ax = c$ . To prove (b), define  $x_\# = x - x_*$ , so that  $x = x_* + x_\#$ . Then since  $Ax = c$  and  $Ax_* = c$ , it follows that

$$Ax_\# = A(x - x_*) = Ax - Ax_* = c - c = 0 \quad \square$$

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  $Ax = 0$  can be expressed as a linear combination of  $r$  linearly independent solutions, where  $r = n - \text{rank}(A)$ .

**Example 6.6.** Consider the system of equations  $Ax = 0$ , where

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

We saw in Example 6.4 that  $A^+A = I_2$ . Thus, the system only has the trivial solution  $0$ .

**Example 6.7.** Since the matrix

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

from Example 6.5 has rank of 2, the homogeneous system of equations  $Ax = 0$  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.

## 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 which 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  $\mathbf{0}$ . One of the most common ways of measuring the closeness of  $Ax_* - c$  to  $\mathbf{0}$  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

$$(Ax_* - c)'(Ax_* - c) \leq (Ax - c)'(Ax - c) \quad (6.4)$$

holds for every  $n \times 1$  vector  $x$ .

Of course, we have already utilized 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 this chapter 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 inverse as the least squares inverse of  $A$ . The following result motivates this description.

**Theorem 6.12.** Let  $A^L$  be any  $\{1, 3\}$ -inverse of a matrix  $A$ . Then the vector  $x_* = A^Lc$  is a least squares solution to the system of equations  $Ax = c$ .

*Proof.* We must show that (6.4) holds when  $x_* = A^Lc$ . The right-hand side of (6.4) can be written as

$$\begin{aligned} (Ax - c)'(Ax - c) &= \{(Ax - AA^Lc) + (AA^Lc - c)\}' \{(Ax - AA^Lc) + (AA^Lc - c)\} \\ &= (Ax - AA^Lc)'(Ax - AA^Lc) + (AA^Lc - c)'(AA^Lc - c) \\ &\quad + 2(Ax - AA^Lc)'(AA^Lc - c) \\ &\geq (AA^Lc - c)'(AA^Lc - c) = (Ax_* - c)'(Ax_* - c), \end{aligned}$$

where the inequality follows from the fact that

$$(Ax - AA^Lc)'(Ax - AA^Lc) \geq 0,$$

and

$$\begin{aligned} (Ax - AA^Lc)'(AA^Lc - c) &= (x - A^Lc)'A'(AA^Lc - c) \\ &= (x - A^Lc)'A'((AA^L)'c - c) \\ &= (x - A^Lc)'(A'A^L'A'c - A'c) \\ &= (x - A^Lc)'(A'c - A'c) = 0 \end{aligned} \quad (6.5)$$

This completes the proof.  $\square$

**Corollary 6.12.1.** The vector  $x_*$  is a least squares solution to the system  $Ax = c$  if and only if

$$(Ax_* - c)'(Ax_* - c) = c'(I_m - AA^L)c$$

*Proof.* From the previous theorem,  $A^Lc$  is a least squares solution for any choice of  $A^L$ , and its sum of squared errors is given by

$$\begin{aligned} (AA^Lc - c)'(AA^Lc - c) &= c'(AA^L - I_m)'(AA^L - I_m)c \\ &= c'(AA^L - I_m)^2c = c'(AA^LAA^L - 2AA^L + I_m)c \\ &= c'(AA^L - 2AA^L + I_m)c = c'(I_m - AA^L)c \end{aligned}$$

The result now follows since, by definition, a least squares solution minimizes the sum of squared errors, and so any other vector  $x_*$  will be a least squares solution if and only if its sum of squared errors is equal to this minimum sum of squares,  $c'(I_m - AA^L)c$ .  $\square$

**Example 6.8.** Let the system of equations  $Ax = c$  have  $A$  and  $c$  given by

$$A = \begin{bmatrix} 1 & 1 & 2 \\ 1 & 0 & 1 \\ 1 & 1 & 2 \\ 2 & 0 & 2 \end{bmatrix}, \quad c = \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix}$$

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

$$AA^Lc = \frac{1}{10} \begin{bmatrix} 5 & 0 & 5 & 0 \\ 0 & 2 & 0 & 4 \\ 5 & 0 & 5 & 0 \\ 0 & 4 & 0 & 8 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 5 \\ 2.2 \\ 5 \\ 4.4 \end{bmatrix} \neq c,$$

it follows from Theorem 6.2 that the system of equations is inconsistent. A least squares solution is then given by

$$A^Lc = \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} = \begin{bmatrix} 2.2 \\ 2.8 \\ 0 \end{bmatrix}$$

Since  $(AA^Lc - c)' = (5, 2.2, 5, 4.4) - (4, 1, 6, 5) = (1, 1.2, -1, -0.6)$ , the sum of squared errors for the least squares solution is

$$(AA^Lc - c)'(AA^Lc - c) = 3.8.$$

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,  $(ABc - c)' = (5, 2.2, 5, 4.4) - (4, 1, 6, 5) = (1, 1.2, -1, -0.6)$ , and so 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 while a least squares solution  $x_*$  may not be unique, the vector  $Ax_*$  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

$$A\mathbf{x}_* = AA^L\mathbf{c} \quad (6.6)$$

*Proof.* Using Theorem 6.2, we see that the system of equations given in (6.6) is consistent since

$$AA^L(AA^L\mathbf{c}) = (AA^LA)A^L\mathbf{c} = AA^L\mathbf{c}$$

The sum of squared errors for any vector  $\mathbf{x}_*$  satisfying (6.6) 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

$$\begin{aligned} (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) &= \mathbf{c}'(I_m - AA^L)\mathbf{c} \\ &= \mathbf{c}'(I_m - AA^L)'(I_m - AA^L)\mathbf{c} \\ &= (AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}), \end{aligned} \quad (6.7)$$

where we have used the fact that  $(I_m - AA^L)$  is symmetric and idempotent. However, we also have

$$\begin{aligned} (A\mathbf{x}_* - \mathbf{c})'(A\mathbf{x}_* - \mathbf{c}) &= \{(A\mathbf{x}_* - AA^L\mathbf{c}) + (AA^L\mathbf{c} - \mathbf{c})\}' \\ &\quad \cdot \{(A\mathbf{x}_* - AA^L\mathbf{c}) + (AA^L\mathbf{c} - \mathbf{c})\} \\ &= (A\mathbf{x}_* - AA^L\mathbf{c})'(A\mathbf{x}_* - AA^L\mathbf{c}) + (AA^L\mathbf{c} - \mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}), \end{aligned} \quad (6.8)$$

since  $(A\mathbf{x}_* - AA^L\mathbf{c})'(AA^L\mathbf{c} - \mathbf{c}) = 0$ , as shown in (6.5). Now (6.7) and (6.8) imply that

$$(A\mathbf{x}_* - AA^L\mathbf{c})'(A\mathbf{x}_* - AA^L\mathbf{c}) = 0,$$

which can be true only if

$$(Ax_* - AA^Lc) = 0,$$

and this establishes (6.6). □

We now give an expression 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

$$x_y = A^Lc + (I_n - A^LA)y,$$

where  $y$  is an arbitrary  $n \times 1$  vector. Then, for each  $y$ ,  $x_y$  is a least squares solution to the system of equations  $Ax = c$ , and for any least squares solution  $x_*$  there exists a vector  $y$  such that  $x_* = x_y$ .

*Proof.* Since

$$A(I_n - A^LA)y = (A - AA^LA)y = (A - A)y = 0,$$

we have  $Ax_y = AA^Lc$ , and so by Theorem 6.13  $x_y$  is a least squares solution. Conversely, if  $x_*$  is an arbitrary least squares solution, then by using Theorem 6.13 again, we must have

$$Ax_* = AA^Lc,$$

which, when premultiplied by  $A^L$ , implies that,

$$0 = -A^LA(x_* - A^Lc)$$

Adding  $x_*$  to both sides of this identity, and then rearranging we get

$$\begin{aligned} x_* &= x_* - A^LA(x_* - A^Lc) \\ &= A^Lc + x_* - A^Lc - A^LA(x_* - A^Lc) \\ &= A^Lc + (I_n - A^LA)(x_* - A^Lc) \end{aligned}$$

This completes the proof since we have shown that  $x_* = x_y$ , where  $y = (x_* - A^Lc)$ . □

We saw in Example 6.8 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  $(I - A^L A) = (0)$ , or equivalently,  $A^L A = I_n$ . The result now follows from Theorem 5.22(g).  $\square$

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 6.16.** 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

$$a'(I_n - A^L A) = 0' \tag{6.9}$$

Now if (6.9) holds, then

$$a' = b'A,$$

where  $b' = a'A^L$ , and so  $a$  is in the row space of  $A$ . On the other hand, if  $a$  is in the row space of  $A$ , then there exists some vector  $b$  such that  $a' = b'A$ . This implies that

$$a'(I_n - A^L A) = b'A(I_n - A^L A) = b'(A - AA^L A) = b'(A - A) = 0',$$

and so the least squares solution must be unique.  $\square$

**Example 6.9.** We will obtain the general least squares solution to the system of equations presented in Example 6.8. First note that

$$A^L A = \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix},$$

so that

$$\begin{aligned} \mathbf{x}_y &= A^L \mathbf{c} + (\mathbf{I}_3 - A^L A) \mathbf{y} \\ &= \frac{1}{10} \begin{bmatrix} 0 & 2 & 0 & 4 \\ 5 & -2 & 5 & -4 \\ 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 4 \\ 1 \\ 6 \\ 5 \end{bmatrix} + \begin{bmatrix} 0 & 0 & -1 \\ 0 & 0 & -1 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_1 \\ y_2 \\ y_3 \end{bmatrix} \\ &= \begin{bmatrix} 2.2 - y_3 \\ 2.8 - y_3 \\ y_3 \end{bmatrix} \end{aligned}$$

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. 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

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad (6.10)$$

where  $\mathbf{y}$  is  $N \times 1$ ,  $X$  is  $N \times m$ ,  $\boldsymbol{\beta}$  is  $m \times 1$ , and  $\boldsymbol{\epsilon}$  is  $N \times 1$ , we have assumed that  $\text{rank}(X) = m$ . In this case, the normal equations,

$$X'X\hat{\boldsymbol{\beta}} = X'\mathbf{y}, \quad (6.11)$$

yield a unique solution, the unique least squares estimator of  $\boldsymbol{\beta}$ , given by

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

However, in many applications, the matrix  $X$  has less than full rank.

**Example 6.10.** Consider the univariate one-way classification model, which was written as



$$y_{ij} = \mu_i + \epsilon_{ij},$$

in Example 3.14, where  $i = 1, \dots, k$  and  $j = 1, \dots, n_i$ . This model can be written in the form of (6.10), where  $\beta = (\mu_1, \dots, \mu_k)'$  and

$$X = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}$$

In this case,  $X$  is of full rank, and so  $\hat{\beta} = (X'X)^{-1}X'y = \bar{y} = (\sum y_{1j}/n_1, \dots, \sum y_{kj}/n_k)'$ . 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  $\mu$  represents an overall effect while  $\tau_i$  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  $y = X_*\beta_* + \epsilon$ , then  $\beta_* = (\mu, \tau_1, \dots, \tau_k)'$  and

$$X_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}$$

Thus, this second parameterization of the one-way classification model has the design matrix  $X_*$  less than full rank since  $\text{rank}(X_*) = k$ .

In this section, we will apply some of the results of this chapter to the estimation of parameters in the model given by (6.10) when  $X$  is less than full rank. First of all, let us consider the task of solving the normal equations given by (6.11); that is, using our usual notation for a system of equations, we want to solve  $Ax = c$ , where  $A = X'X$ ,  $x = \hat{\beta}$ , and  $c = X'y$ . Now from Theorem 6.2, we see that (6.11) is a consistent system of equations since

$$\begin{aligned} X'X(X'X)^+X'y &= X'XX^+X'^+X'y = X'XX^+(XX^+)'y = X'XX^+XX^+y \\ &= X'XX^+y = X'(XX^+)'y = X'X'^+X'y = X'y \end{aligned}$$

Consequently, using Theorem 6.4, we find that the general solution  $\hat{\beta}$  can be written as

$$\hat{\beta} = (X'X)^- X'y + \{I - (X'X)^- X'X\}u, \quad (6.12)$$

or, if we use the Moore–Penrose generalized inverse, as

$$\begin{aligned} \hat{\beta} &= (X'X)^+ X'y + \{I - (X'X)^+ X'X\}u \\ &= X^+ y + (I - X^+ X)u, \end{aligned}$$

where  $u$  is an arbitrary  $m \times 1$  vector. The same general solution can be obtained by using the least squares results of Section 6.5 on the system of equations

$$y = X\hat{\beta}$$

Thus, using Theorem 6.14 with  $A = X$ ,  $x = \hat{\beta}$ , and  $c = y$ , the least squares solution is given by

$$\hat{\beta} = X^L y + (I - X^L X)u,$$

which is, of course, equivalent to that given by (6.12).

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  $y = X\beta + \epsilon$  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  $a'\beta$  of the parameter vector  $\beta$  is estimable if and only if there exists some  $N \times 1$  vector  $b$  such that

$$a'\beta = E(b'y) = b'E(y) = b'X\beta;$$

that is, if and only if there exists a linear function of the components of  $y$ ,  $b'y$ , which is an unbiased estimator of  $a'\beta$ .

The condition that a linear function  $a'\beta$  be estimable is equivalent to the condition that the corresponding estimator  $a'\hat{\beta}$  be unique. To see this, note that from the definition above, the function  $a'\beta$  is estimable if and only if  $a$  is in the row space of  $X$ , while it follows from Theorem 6.16 that  $a'\hat{\beta}$  is unique

if and only if  $a$  is in the row space of  $X$ . In addition, since  $X'(XX')^+X$  is the projection matrix for the row space of  $X$ , we get the more practical condition for estimability of  $a'\beta$  given by

$$X'(XX')^+Xa = a \tag{6.13}$$

It follows from Theorems 5.3 and 5.25 that

$$X'(XX')^+X = X'X'^+ = X'X'^L = X'(XX')^-X,$$

and so equation (6.13) is not dependent upon 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  $\hat{y} = X\hat{\beta}$  and its sum of squared errors  $(y - \hat{y})'(y - \hat{y})$  to the choice of the least squares solution  $\hat{\beta}$ . Since  $XX^+X = X$

$$\hat{y} = X\hat{\beta} = X\{X^+y + (I - X^+X)u\} = XX^+y + (X - XX^+X)u = XX^+y,$$

which does not depend on the vector  $u$ . Thus,  $\hat{y}$  is unique, while the uniqueness of

$$(y - \hat{y})'(y - \hat{y}) = y'(I - XX^+)y$$

follows immediately from the uniqueness of  $\hat{y}$ .

**Example 6.11.** Let us return to the one-way classification model

$$y = X_*\beta_* + \epsilon$$

of Example 6.10, where  $\beta_* = (\mu, \tau_1, \dots, \tau_k)'$  and

$$X_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix}$$

Since the rank of the  $n \times (k + 1)$  matrix  $X_*$ , where  $n = \sum n_i$  is  $k$ , the least squares solution for  $\beta_*$  is not unique. To find the form of the general solution, note that

$$X'_*X_* = \begin{bmatrix} n & n' \\ n & D_n \end{bmatrix},$$

while 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  $n = (n_1, \dots, n_k)'$  and  $D_n = \text{diag}(n_1, \dots, n_k)$ . Thus, using (6.12) we have the general solution

$$\begin{aligned} \hat{\beta}_* &= \begin{bmatrix} n^{-1} & \mathbf{0}' \\ \mathbf{0} & D_n^{-1} - n^{-1}\mathbf{1}_k\mathbf{1}_k' \end{bmatrix} \begin{bmatrix} n\bar{y} \\ D_n\bar{y} \end{bmatrix} + \left\{ \mathbf{1}_{k+1} - \begin{bmatrix} 1 & n^{-1}n' \\ \mathbf{0} & \mathbf{I}_k - n^{-1}\mathbf{1}_k n' \end{bmatrix} \right\} u \\ &= \begin{bmatrix} \bar{y} \\ \bar{y} - \bar{y}\mathbf{1}_k \end{bmatrix} + \begin{bmatrix} 0 & -n^{-1}n' \\ \mathbf{0} & n^{-1}\mathbf{1}_k n' \end{bmatrix} u, \end{aligned}$$

where  $\bar{y} = (\bar{y}_1, \dots, \bar{y}_k)'$  and  $\bar{y} = \sum n_i \bar{y}_i / n$ . Choosing  $u = \mathbf{0}$ , 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, since  $\mu + \tau_i = \mathbf{a}'_i\beta_*$ , where  $\mathbf{a}_i = (1, \mathbf{e}'_i)'$ , its estimator is given by

$$\mathbf{a}'_i\hat{\beta}_* = [1 \quad \mathbf{e}'_i] \begin{bmatrix} \bar{y} \\ \bar{y} - \bar{y} \end{bmatrix} = \bar{y}_i$$

The vector of fitted values is

$$\hat{y} = X_*\hat{\beta}_* = \begin{bmatrix} \mathbf{1}_{n_1} & \mathbf{1}_{n_1} & \mathbf{0} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_2} & \mathbf{0} & \mathbf{1}_{n_2} & \cdots & \mathbf{0} \\ \mathbf{1}_{n_3} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{0} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \mathbf{1}_{n_k} & \mathbf{0} & \mathbf{0} & \cdots & \mathbf{1}_{n_k} \end{bmatrix} \begin{bmatrix} \bar{y} \\ \bar{y}_1 - \bar{y} \\ \bar{y}_2 - \bar{y} \\ \vdots \\ \bar{y}_k - \bar{y} \end{bmatrix} = \begin{bmatrix} \bar{y}_1 \mathbf{1}_{n_1} \\ \bar{y}_2 \mathbf{1}_{n_2} \\ \vdots \\ \bar{y}_k \mathbf{1}_{n_k} \end{bmatrix},$$

while the sum of squared errors is given by

$$(\mathbf{y} - \hat{\mathbf{y}})'(\mathbf{y} - \hat{\mathbf{y}}) = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)^2$$

### 7. SYSTEMS OF LINEAR EQUATIONS AND THE SINGULAR VALUE DECOMPOSITION

When  $A$  is square and nonsingular, then the solution to the system of equations  $Ax = c$  can be conveniently expressed in terms of the inverse of  $A$ , as  $x = A^{-1}c$ . 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

$$Dy = b, \tag{6.14}$$

where  $y$  is an  $n \times 1$  vector of variables,  $b$  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  $\Delta$  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 (6.14) is consistent with the unique solution given by  $y = \Delta^{-1}b$ . For (b), if we partition  $y$  as  $y = (y_1', y_2')'$ , where  $y_1$  is  $r \times 1$ , then equation (6.14) reduces to

$$\Delta y_1 = b$$

Thus, (6.14) is consistent and has solutions of the form

$$y = \begin{bmatrix} \Delta^{-1}b \\ y_2 \end{bmatrix},$$

where the  $(n - r) \times 1$  vector  $y_2$  is arbitrary. Since we then have  $n - r$  linearly independent choices for  $y_2$ , the number of linearly independent solutions is  $n - r$  if  $b = 0$  and  $n - r + 1$  if  $b \neq 0$ . When  $D$  has the form given in (c), the system in (6.14) takes the form

$$\begin{bmatrix} \Delta y \\ 0 \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix},$$

where  $b_1$  is  $r \times 1$  and  $b_2$  is  $(m - r) \times 1$ , and so it is consistent only if  $b_2 = 0$ . If

this is the case, the system then has a unique solution given by  $y = \Delta^{-1}b_1$ . For the final form given in (d), the system of equations in (6.14) appears as

$$\begin{bmatrix} \Delta y_1 \\ \mathbf{0} \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix},$$

where  $y$  and  $b$  have been partitioned as before. As in the case of form (c), this system is consistent only if  $b_2 = \mathbf{0}$ , and as in the case of form (b), when consistent, it has  $n - r$  linearly independent solutions if  $b = \mathbf{0}$ , and  $n - r + 1$  linearly independent solutions if  $b \neq \mathbf{0}$ . The general solution is given by

$$y = \begin{bmatrix} \Delta^{-1}b_1 \\ y_2 \end{bmatrix},$$

where the  $(n - r) \times 1$  vector  $y_2$  is arbitrary.

All of the above can now be readily applied to the general system of equations,

$$Ax = c \tag{6.15}$$

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 (6.14), where the vector of variables is given by  $y = Q'x$  and the vector of constants is  $b = P'c$ . Consequently, if  $y$  is a solution to (6.14), then  $x = Qy$  will be a solution to (6.15). Thus, in the case of forms (a) and (b), (6.15) is consistent with the unique solution given by

$$x = Qy = Q\Delta^{-1}b = Q\Delta^{-1}P'c = A^{-1}c,$$

when (a) is the form of  $D$ , while for form (b) the general solution is

$$x = Qy = [Q_1 \quad Q_2] \begin{bmatrix} \Delta^{-1}b \\ y_2 \end{bmatrix} = Q_1\Delta^{-1}P'c + Q_2y_2,$$

where  $Q_1$  is  $n \times r$  and  $y_2$  is an arbitrary  $(n - r) \times 1$  vector. The term  $Q_2y_2$  has no effect on the value of  $Ax$  since 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 (6.15) is consistent only if  $c = P_1b_1$  so that  $P_2b_2 = \mathbf{0}$ , where  $P = (P_1, P_2)$  and  $P_1$  is  $m \times r$ ; that is, since the columns of  $P_1$  form a basis for the range of  $A$ , the system is consistent if  $c$  is in the column space of  $A$ . Thus, if we partition  $c$  as  $c = (c'_1, c'_2)'$  where  $c_1$  is  $r \times 1$ , then when form (c) holds, the unique solution will be given by

$$x = Qy = Q\Delta^{-1}b_1 = Q\Delta^{-1}P_1'c$$

In the case of form (d), the general solution is

$$x = Qy = [Q_1 \quad Q_2] \begin{bmatrix} \Delta^{-1}b_1 \\ y_2 \end{bmatrix} = Q_1\Delta^{-1}P_1'c + Q_2y_2$$

### 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, utilizes 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 since we only need store the nonzero values and their location. 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, Erisman, and Reid (1986) and Golub and Van Loan (1989).

A second approach to the solution of sparse systems of equations utilizes iterative methods. In this case, a sequence of vectors,  $x_0, x_1, \dots$  is generated with  $x_0$  being some initial vector while  $x_j$  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  $Ax = c$ . 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  $Ax = c$  can be written as

$$(A - D_A)x + D_A x = c,$$

which yields the identity

$$x = D_A^{-1} \{c - (A - D_A)x\}$$

This is the motivation for the Jacobi method that computes  $x_j$  as

$$x_j = D_A^{-1} \{c - (A - D_A)x_{j-1}\}$$

On the other hand, the Gauss-Seidel method utilizes 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,  $Ax = c$  can be rearranged as

$$A_1 x = c - A_2 x,$$

and this leads to the iterative scheme

$$A_1 x_j = c - A_2 x_{j-1},$$

which is easily solved for  $x_j$  since 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 (1989).

Consider the function

$$f(x) = \frac{1}{2}x'Ax - x'c,$$

where  $x$  is an  $m \times 1$  vector and  $A$  is an  $m \times m$  positive definite matrix. The vector of partial derivatives of  $f(x)$  given by

$$\nabla f(x) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_m} \right)' = Ax - c$$

is sometimes referred to as the gradient of  $f(x)$ . Setting this equal to the zero vector, we find that the vector minimizing  $f$ ,  $x = A^{-1}c$ , is also the solution to the system  $Ax = c$ . Thus, a vector which approximately minimizes  $f$  will also be an approximate solution to  $Ax = c$ . One iterative method for finding the minimizer  $x$  involves successively finding minimizers  $x_j$  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,  $q_1, \dots, q_m$ , we will define the  $j$ th subspace as the space with the columns of the  $m \times j$  matrix,  $Q_j = (q_1, \dots, q_j)$ , as its basis. Consequently, for some  $j \times 1$  vector  $y_j$ ,

$$x_j = Q_j y_j \tag{6.16}$$

and

$$f(x_j) = \min_{y \in R^j} f(Q_j y) = \min_{y \in R^j} g(y) = g(y_j),$$

where

$$g(y) = \frac{1}{2}y'(Q_j' A Q_j)y - y' Q_j' c$$

Thus, the gradient of  $g(y_j)$  must be equal to the null vector, and so

$$(Q_j' A Q_j)y_j = Q_j' c \tag{6.17}$$

To obtain  $x_j$ , we can first use (6.17) to calculate  $y_j$  and then use this in (6.16) to

get  $x_j$ . The final  $x_j$ ,  $x_m$ , will be the solution to  $Ax = c$ , but the goal here is to stop the iterative process before  $j = m$  with a sufficiently accurate solution  $x_j$ .

The iterative scheme described above will work with different sets of orthonormal vectors  $q_1, \dots, q_m$  but we will see that by a judicious choice of this set, we may guarantee that the computation involved in computing the  $x_j$ 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  $Ax = c$ . Let  $\lambda_1$  and  $\lambda_m$  denote the largest and smallest eigenvalues of  $A$ , while  $\lambda_{1j}$  and  $\lambda_{jj}$  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  $\lambda_1$  and  $\lambda_m$  are the maximum and minimum values of the Rayleigh quotient,

$$R(x, A) = \frac{x'Ax}{x'x}$$

Suppose that we have the  $j$  columns of  $Q_j$ , and we wish to find an additional vector  $q_{j+1}$  so as to form the matrix  $Q_{j+1}$  and have  $\lambda_{1,j+1}$  and  $\lambda_{j+1,j+1}$  as close to  $\lambda_1$  and  $\lambda_m$  as possible. If  $u_j$  is a vector in the space spanned by the columns of  $Q_j$  and satisfying  $R(u_j, A) = \lambda_{1j}$ , then since the gradient

$$\nabla R(u_j, A) = \frac{2}{u_j' u_j} \{ Au_j - R(u_j, A)u_j \}$$

gives the direction in which  $R(u_j, A)$  is increasing most rapidly, we would want to choose  $q_{j+1}$  so that  $\nabla R(u_j, A)$  is in the space spanned by the columns of  $Q_{j+1}$ . On the other hand, if  $v_j$  is a vector in the space spanned by  $Q_j$  and satisfying  $R(v_j, A) = \lambda_{jj}$ , then since  $R(v_j, A)$  is decreasing most rapidly in the direction given by  $-\nabla R(v_j, A)$ , we would want to make sure that  $\nabla R(v_j, A)$  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  $q_1, Aq_1, \dots, A^{j-1}q_1$  and we select  $q_{j+1}$  so that the columns of  $Q_{j+1}$  are spanned by the vectors  $q_1, Aq_1, \dots, A^j q_1$ , since both  $\nabla R(u_j, A)$  and  $\nabla R(v_j, A)$  are of the form  $aAx + bx$  for some vector  $x$  spanned by the columns of  $Q_j$ . Thus, we start with an initial unit vector  $q_1$ , while for  $j \geq 2$ ,  $q_j$  is selected as a unit vector orthogonal to  $q_1, \dots, q_{j-1}$  and such that the columns of  $Q_j$  are spanned by the vectors  $q_1, Aq_1, \dots, A^{j-1}q_1$ . These particular  $q_j$  vectors are known as the Lanczos vectors. The calculation of the  $q_j$ 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 choose  $P$  and  $q_1$  so that  $Pe_1 = q_1$ , then

$$(q_1, Aq_1, \dots, A^{j-1}q_1) = P(e_1, Te_1, \dots, T^{j-1}e_1).$$

Since  $(e_1, Te_1, \dots, T^{j-1}e_1)$  has upper triangular structure, this means that the first  $j$  columns of  $P$  span the column space of  $(q_1, Aq_1, \dots, A^{j-1}q_1)$ ; that is, the  $q_j$ s can be obtained by calculating the factorization  $A = PTP'$ , or in other words, we can take  $Q = (q_1, \dots, q_m) = P$ . Thus, since  $AQ = QT$ , we have

$$Aq_1 = \alpha_1q_1 + \beta_1q_2, \tag{6.18}$$

and

$$Aq_j = \beta_{j-1}q_{j-1} + \alpha_jq_j + \beta_jq_{j+1}, \tag{6.19}$$

for  $j = 2, \dots, m - 1$ . Using these equations and the orthonormality of the  $q_j$ s, it is easily shown that  $\alpha_j = q_j' Aq_j$  for all  $j$ , and as long as  $p_j = (A - \alpha_j I_m)q_j - \beta_{j-1}q_{j-1} \neq 0$ , then  $\beta_j^2 = p_j' p_j$  and  $q_{j+1} = p_j / \beta_j$  for  $j = 1, \dots, m - 1$ , if we define  $q_0 = 0$ . Thus, we can continue calculating the  $q_j$ s until we encounter a  $p_j = 0$ . To see the significance of this event, let us suppose that the iterative procedure has proceeded through the first  $j - 1$  steps with  $p_i \neq 0$  for each  $i = 2, \dots, j - 1$ , and so we have obtained the matrix  $Q_j$  whose columns form a basis for  $(q_1, Aq_1, \dots, A^{j-1}q_1)$ . Note that it follows immediately from the relationship  $AQ = QT$  that

$$AQ_j = Q_j T_j + p_j e_j',$$

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  $Q_j' A Q_j = T_j + Q_j' p_j e_j'$ . But  $q_i' A q_i = \alpha_i$ , while it follows from (6.18) and (6.19) that  $q_{i+1}' A q_i = \beta_i$  and  $q_k' A q_i = 0$  if  $k > i + 1$ . Thus,  $Q_j' A Q_j = T_j$  and so we must have  $Q_j' p_j = 0$ . Now if  $p_j \neq 0$ , then  $q_{j+1} = p_j / \beta_j$  is orthogonal to the columns of  $Q_j$ . Further, it follows from the fact that  $q_{j+1}$  is a linear combination of  $Aq_j, q_j$ , and  $q_{j-1}$  that the columns of  $Q_{j+1} = (Q_j, q_{j+1})$  form a basis for the column space of  $(q_1, Aq_1, \dots, A^j q_1)$ . If, on the other hand,  $p_j = 0$ , then  $AQ_j = Q_j T_j$ . From this we see that the vectors  $A^j q_1, \dots, A^{m-1} q_1$  are in the space spanned by the columns of  $Q_j$ , that is, the space spanned by the

vectors  $q_1, Aq_1, \dots, A^{j-1}q_1$ . Consequently, the iterative procedure is complete since there are only  $j$   $q_i$ 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 due to the encounter of a  $p_j = 0$ , but due to sufficiently accurate approximations of the eigenvalues of  $A$ .

Now let us return to the problem of solving the system of equations  $Ax = c$  through the iterative procedure based on the calculation of  $y_j$  in (6.17) and then  $x_j$  in (6.16). 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 (6.17) 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  $y_j$  in (6.17) can be easily found by first solving  $L_j w_j = Q_j' c$ , then  $D_j z_j = w_j$ , and finally  $L_j' y_j = z_j$ . Even as  $j$  increases, the computation required is not extensive since  $D_{j-1}$  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  $x_j$  from  $y_j$  using (6.16). 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  $B_j = (b_1, \dots, b_j)$  so that  $B_j L_j' = Q_j$ , then by premultiplying the equation  $T_j y_j = Q_j' c$  by  $Q_j T_j^{-1}$  and using (6.16), we get

$$x_j = Q_j T_j^{-1} Q_j' c = Q_j (L_j D_j L_j')^{-1} Q_j' c = B_j z_j, \quad (6.20)$$

where  $z_j$  is as previously defined. It will be easier to compute  $x_j$  from (6.20) than from (6.16) since  $B_j$  and  $z_j$  are simple to compute after  $B_{j-1}$  and  $z_{j-1}$  have already been calculated. For instance, from the definition of  $B_j$ , we see that  $b_1 = q_1$  and  $b_i = q_i - l_{i-1} b_{i-1}$  for  $i > 1$ , and consequently,  $B_j = (B_{j-1}, b_j)$ . Using the defining equations for  $w_j$  and  $z_j$ , we find that

$$L_j D_j z_j = Q_j' c \quad (6.21)$$

## PROBLEMS

If we partition  $z_j$  as  $z_j = (\gamma'_{j-1}, \gamma_j)'$ , where  $\gamma_{j-1}$  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 (6.21) implies that  $L_{j-1}D_{j-1}\gamma_{j-1} = Q'_{j-1}c$ . But this means that  $\gamma_{j-1} = z_{j-1}$ , and so to compute  $z_j$ , we only need to compute  $\gamma_j$ , which is given by

$$\gamma_j = (q'_j c - l_{j-1} d_{j-1} \gamma_{j-1}) / d_j,$$

where  $\gamma_{j-1}$  is the last component of  $z_{j-1}$ . Thus, (6.20) becomes

$$x_j = B_j z_j = [B_{j-1}, b_j] \begin{bmatrix} z_{j-1} \\ \gamma_j \end{bmatrix} = B_{j-1} z_{j-1} + \gamma_j b_j = x_{j-1} + \gamma_j b_j,$$

and so we have a simple formula for computing the  $j$ th iterative solution from  $b_j$ ,  $\gamma_j$ , and the  $(j-1)$ th iterative solution  $x_{j-1}$ .

## PROBLEMS

1. Consider the system of equations  $Ax = c$ , where  $A$  is the  $4 \times 3$  matrix given in Problem 5.2 and

$$c = \begin{bmatrix} 1 \\ 3 \\ -1 \\ 0 \end{bmatrix}$$

- (a) Show that the system is consistent.
- (b) Find a solution to this system of equations.
- (c) How many linearly independent solutions are there?

2. The system of equations  $Ax = c$  has  $A$  equal to the  $3 \times 4$  matrix given in Problem 5.36 and

$$c = \begin{bmatrix} 1 \\ 1 \\ 4 \end{bmatrix}$$

- (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.

3. Suppose the system of equations  $Ax = c$  has

$$A = \begin{bmatrix} 5 & 2 & 1 \\ 3 & 1 & 1 \\ 2 & 1 & 0 \\ 1 & 2 & -3 \end{bmatrix}$$

For each  $c$  given below, determine whether or not the system of equations is consistent.

$$(a) \ c = \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix}, \quad (b) \ c = \begin{bmatrix} 3 \\ 2 \\ 1 \\ -1 \end{bmatrix}, \quad (c) \ c = \begin{bmatrix} 1 \\ -1 \\ 1 \\ -1 \end{bmatrix}$$

4. Consider the system of equations  $Ax = c$ , where

$$A = \begin{bmatrix} 1 & 1 & -1 & 0 & 2 \\ 2 & 1 & 1 & 1 & 1 \end{bmatrix}, \quad c = \begin{bmatrix} 3 \\ 1 \end{bmatrix}$$

- (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.

5. Prove Theorem 6.5.

6. Consider the system of equations  $AXB = C$ , where  $X$  is a  $3 \times 3$  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.

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  $c \neq \mathbf{0}$  and  $y \neq \mathbf{0}$ .
8. Suppose the  $m \times n$  matrix  $A$  and  $m \times 1$  vector  $c \neq \mathbf{0}$  are such that  $A^{-}c$  is the same for all choices of  $A^{-}$ . Use Theorem 5.23 to show that, if  $Ax = c$  is a consistent system of equations, then it has a unique solution.
9. For the homogeneous system of equations  $Ax = \mathbf{0}$  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.

10. 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.
11. Let

$$A = \begin{bmatrix} 1 & -1 & 1 & 1 \\ 2 & 3 & 1 & -1 \end{bmatrix}, \quad c = \begin{bmatrix} 1 \\ 2 \end{bmatrix},$$

$$B = \begin{bmatrix} 2 & 1 & 2 & -1 \\ 0 & 1 & 1 & 1 \end{bmatrix}, \quad d = \begin{bmatrix} 2 \\ 4 \end{bmatrix}.$$

- (a) Show that the system  $Ax = c$  is consistent and has three linearly independent solutions.
- (b) Show that the system  $Bx = d$  is consistent and has three linearly independent solutions.
- (c) Show that the systems  $Ax = c$  and  $Bx = d$  have a common solution and that this common solution is unique.
12. 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

$$X_* = A^{-}C + (I - A^{-}A)DB^{-} + (I - A^{-}A)Y(I - BB^{-}),$$

where  $Y$  is an arbitrary  $n \times p$  matrix.

13. In Exercise 5.37, a least squares inverse was found for the matrix

$$A = \begin{bmatrix} 1 & -1 & -2 & 1 \\ -2 & 4 & 3 & -2 \\ 1 & 1 & -3 & 1 \end{bmatrix}$$

- Use this least squares inverse to show that the system of equations  $Ax = c$  is inconsistent, where  $c' = (2, 1, 5)$ .
- Find a least squares solution.
- Compute the sum of squared errors for a least squares solution to this system of equations.

14. Consider the system of equations  $Ax = c$ , where

$$A = \begin{bmatrix} 1 & 0 & 2 \\ 2 & -1 & 3 \\ -1 & 2 & 0 \\ -2 & 1 & -3 \end{bmatrix}, \quad c = \begin{bmatrix} 2 \\ 2 \\ 5 \\ 0 \end{bmatrix}$$

- Find a least squares inverse of  $A$ .
  - Show that the system of equations is inconsistent.
  - Find a least squares solution.
  - Is this solution unique?
15. Show that  $x_*$  is a least squares solution to the system of equations  $Ax = c$  if and only if

$$A'Ax_* = A'c$$

16. Let  $A$  be an  $m \times n$  matrix, and  $x_*$ ,  $y_*$ , and  $c$ ,  $n \times 1$ ,  $m \times 1$ , and  $m \times 1$  vectors, respectively. Suppose that  $x_*$  and  $y_*$  are such that the system of equations

$$\begin{bmatrix} I_m & A \\ A' & (0) \end{bmatrix} \begin{bmatrix} y_* \\ x_* \end{bmatrix} = \begin{bmatrix} c \\ 0 \end{bmatrix},$$

holds. Show that  $x_*$  then must be a least squares solution to the system  $Ax = c$ .

17. 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  $\mu$  represents an overall effect,  $\tau_i$  is an effect due to the  $i$ th level of factor one,  $\gamma_j$  is an effect due to the  $j$ th level of factor two, and  $\eta_{ij}$  is an effect due to the interaction of the  $i$ th and  $j$ th levels of factors one and two; as usual, the  $e_{ijk}$ s represent independent random errors, each distributed as  $N(0, \sigma^2)$ .

- (a) Set up the vectors  $y$ ,  $\beta$ , and  $\epsilon$  and the matrix  $X$  so that the two-way model above can be written in the matrix form  $y = X\beta + \epsilon$ .
- (b) Find the rank  $r$ , of  $X$ . Determine a set of  $r$  linear independent estimable functions of the parameters,  $\mu$ ,  $\tau_i$ ,  $\gamma_j$ , and  $\eta_{ij}$ .
- (c) Find a least squares solution for the parameter vector  $\beta$ .

18. Consider the regression model

$$y = X\beta + \epsilon,$$

where  $X$  is  $N \times m$ ,  $\epsilon \sim N_N(\mathbf{0}, \sigma^2 C)$ , 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  $\hat{\beta} = (X' C^{-1} X)^{-1} X' C^{-1} y$ , which minimizes

$$(y - X\hat{\beta})' C^{-1} (y - X\hat{\beta}) \tag{6.22}$$

Show that if  $X$  is less than full column rank, then the generalized least squares estimator of  $\beta$  which minimizes (6.22) is given by

$$\hat{\beta} = (X' C^{-1} X)^- X' C^{-1} y + \{I - (X' C^{-1} X)^- X' C^{-1} X\} u,$$

where  $u$  is an arbitrary  $m \times 1$  vector.

19. Restricted least squares obtains the vectors  $\hat{\beta}$  that minimize

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

subject to the restriction that  $\hat{\beta}$  satisfies  $B\hat{\beta} = b$ , where  $B$  is  $p \times m$  and  $b$  is  $p \times 1$  such that  $BB^-b = b$ . Use Theorem 6.4 to find the general solution  $\hat{\beta}_u$  to the consistent system of equations  $B\hat{\beta} = b$ , where  $\hat{\beta}_u$  depends on an arbitrary vector  $u$ . Substitute this expression for  $\hat{\beta}$  into  $(y - X\hat{\beta})'(y - X\hat{\beta})$ , and then use Theorem 6.14 to obtain the general least squares solution  $u_w$ , for  $u$ , where  $u_w$  depends on an arbitrary vector  $w$ . Substitute  $u_w$  for  $u$  in  $\hat{\beta}_u$  to show that the general restricted least squares solution for  $\beta$  is given by

$$\begin{aligned} \hat{\beta}_w = & B^-b + (I - B^-B)\{[X(I - B^-B)]^L(y - XB^-b) \\ & + (I - [X(I - B^-B)]^L X(I - B^-B))w\}. \end{aligned}$$

20. In the previous exercise, show that if we use the Moore–Penrose inverse as the least squares inverse of  $[X(I - B^{-1}B)]$  in the expression given for  $\hat{\beta}_w$ , then it simplifies to

$$\hat{\beta}_w = B^{-1}b + [X(I - B^{-1}B)]^+(y - XB^{-1}b) + (I - B^{-1}B)\{I - [X(I - B^{-1}B)]^+X(I - B^{-1}B)\}w.$$

21. Consider the iterative procedure, based on the Lanczos vectors, for solving the system of equations  $Ax = c$ . Suppose that for the initial Lanczos vector  $q_1$  we use  $(c'c)^{-1/2}c$ .

- (a) Show that if for some  $j$ ,  $p_j = (A - \alpha_j I_m)q_j - \beta_{j-1}q_{j-1} = 0$ , then  $Ax_j = c$ .  
 (b) Show that for any  $j$ , the procedure easily yields a measure of the adequacy of the  $j$ th iterative solution since

$$(Ax_j - c)'(Ax_j - c) = \beta_j^2 y_{jj}^2,$$

where  $y_{jj}$  is the  $j$ th component of the vector  $y_j$  in (6.16).

## CHAPTER SEVEN

# Special Matrices and Matrix Operators

### 1. INTRODUCTION

The concept of partitioning matrices was first introduced in Chapter 1, and we have subsequently used partitioned matrices throughout this text. In this chapter we develop some specialized formulas for the determinant and inverse of partitioned matrices. In addition to partitioned matrices, we will look at some other special types of structured matrices that we have not previously discussed. In this chapter, we will also introduce and develop properties of some special matrix operators. In many situations, a seemingly complicated matrix expression can be written in a fairly simple form by making use of one or more of these matrix operators.

### 2. PARTITIONED MATRICES

Up to this point, most of our applications involving partitioned matrices have utilized only the simple operations of matrix addition and multiplication. In this section, we will obtain expressions for the inverse and determinant of an  $m \times m$  matrix  $A$  that is partitioned into the  $2 \times 2$  block form given by

$$A = \begin{bmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{bmatrix}, \quad (7.1)$$

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  $m_2 \times m_2$ . We wish to obtain expressions for the inverse and determinant of  $A$  in terms of its submatrices. We begin with the inverse of  $A$ .

**Theorem 7.1.** Let the  $m \times m$  matrix  $A$  be partitioned as in (7.1), and suppose that  $A$ ,  $A_{11}$ , and  $A_{22}$  are nonsingular matrices. 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 we have

- (a)  $B_{11} = (A_{11} - A_{12}A_{22}^{-1}A_{21})^{-1} = A_{11}^{-1} + A_{11}^{-1}A_{12}B_{22}A_{21}A_{11}^{-1}$ ,  
 (b)  $B_{22} = (A_{22} - A_{21}A_{11}^{-1}A_{12})^{-1} = A_{22}^{-1} + A_{22}^{-1}A_{21}B_{11}A_{12}A_{22}^{-1}$ ,  
 (c)  $B_{12} = -A_{11}^{-1}A_{12}B_{22}$ ,  
 (d)  $B_{21} = -A_{22}^{-1}A_{21}B_{11}$ .

*Proof.* 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) and (7.5) for  $B_{12}$  and  $B_{21}$ , respectively, immediately leads to the expressions given in (c) and (d). Substituting these solutions for  $B_{12}$  and  $B_{21}$  into (7.2) and (7.3) and solving for  $B_{11}$  and  $B_{22}$  yields the first expressions given for  $B_{11}$  and  $B_{22}$  in (a) and (b). The second expressions in (a) and (b) follow immediately from the first after using Theorem 1.7.  $\square$

**Example 7.1.** Consider the regression model

$$y = X\beta + \epsilon,$$

where  $y$  is  $N \times 1$ ,  $X$  is  $N \times (k+1)$ ,  $\beta$  is  $(k+1) \times 1$ , and  $\epsilon$  is  $N \times 1$ . Suppose that  $\beta$  and  $X$  are partitioned as  $\beta = (\beta'_1, \beta'_2)'$  and  $X = (X_1, X_2)$  so that the product  $X_1\beta_1$  is defined, and we are interested in comparing the complete regression model given above to the reduced regression model

$$y = X_1\beta_1 + \epsilon$$

If  $X$  has full column rank, then the least squares estimators for the two models

are  $\hat{\beta} = (X'X)^{-1}X'y$  and  $\hat{\beta}_1 = (X_1'X_1)^{-1}X_1'y$ , respectively, and the difference in the sums of squared errors for the two models

$$\begin{aligned} & (y - X_1\hat{\beta}_1)'(y - X_1\hat{\beta}_1) - (y - X\hat{\beta})'(y - X\hat{\beta}) \\ &= y'(I - X_1(X_1'X_1)^{-1}X_1')y - y'(I - X(X'X)^{-1}X')y \\ &= y'X(X'X)^{-1}X'y - y'X_1(X_1'X_1)^{-1}X_1'y \end{aligned} \quad (7.6)$$

gives the reduction in the sum of squared errors attributable to the inclusion of the term  $X_2\beta_2$  in the complete model. By using the geometrical properties of least squares regression in Example 2.10, we showed that this reduction in the sum of squared errors simplifies to

$$y'X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}'y,$$

where  $X_{2*} = (I - X_1(X_1'X_1)^{-1}X_1')X_2$ . 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 (7.6) and then simplifying, we get  $y'X_{2*}(X_{2*}'X_{2*})^{-1}X_{2*}'y$ , as required.

Before obtaining expressions for the determinant of  $A$ , 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 use 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 the following.

**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  $A$  in the general case.

**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, and
- (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 similar.  $\square$

**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 3.18(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 3.18(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, since  $A$  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  $c = (ab - m^2)^{-1}$ .

We will use Theorem 7.4 to establish the following useful result.

**Theorem 7.5.** 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

$$\begin{vmatrix} I_m & A \\ -B & I_n \end{vmatrix} = |I_m + AB| \quad (7.7)$$

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

$$\begin{vmatrix} I_m & A \\ -B & I_n \end{vmatrix} = |I_n + BA| \quad (7.8)$$

The result now follows by equating (7.7) and (7.8).  $\square$

Our final result follows directly from Theorem 7.5 if we replace  $A$  by  $-\lambda A$ .



**Corollary 7.5.1.** 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  $BA$ .

### 3. 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

$$\begin{bmatrix} a_{11}B & a_{12}B & \cdots & a_{1n}B \\ a_{21}B & a_{22}B & \cdots & a_{2n}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1}B & a_{m2}B & \cdots & a_{mn}B \end{bmatrix} \quad (7.9)$$

The Kronecker product defined above is more precisely known as the right Kronecker product, and it is the most common definition of 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 (7.9). Throughout this book, any reference to the Kronecker product will be referring to the right Kronecker product. The special structure of the matrix given in (7.9) 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 7.3.** Let  $A$  and  $B$  be the  $1 \times 3$  and  $2 \times 2$  matrices given by

$$A = [0 \quad 1 \quad 2], \quad B = \begin{bmatrix} 1 & 2 \\ 3 & 4 \end{bmatrix}$$

Then we find that

$$A \otimes B = [0B \quad 1B \quad 2B] = \begin{bmatrix} 0 & 0 & 1 & 2 & 2 & 4 \\ 0 & 0 & 3 & 4 & 6 & 8 \end{bmatrix},$$

while

$$B \otimes A = \begin{bmatrix} 1A & 2A \\ 3A & 4A \end{bmatrix} = \begin{bmatrix} 0 & 1 & 2 & 0 & 2 & 4 \\ 0 & 3 & 6 & 0 & 4 & 8 \end{bmatrix}$$

Some of the basic properties of the Kronecker product, which are easily proven from its definition, are summarized below. The proofs are left to the reader as an exercise.

**Theorem 7.6.** Let  $A$ ,  $B$ , and  $C$  be any matrices and  $a$  and  $b$  be any two vectors. Then

- (a)  $\alpha \otimes A = A \otimes \alpha = \alpha A$ , for any scalar  $\alpha$ ,
- (b)  $(\alpha A) \otimes (\beta B) = \alpha\beta(A \otimes B)$ , for any scalars  $\alpha$  and  $\beta$ ,
- (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)  $ab' = a \otimes b' = b' \otimes a$ .

We have the following very useful property involving the Kronecker product and ordinary matrix multiplication.

**Theorem 7.7.** 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

$$(A \otimes B)(C \otimes D) = AC \otimes BD \quad (7.10)$$

*Proof.* The left-hand side of (7.10) 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

$$F_{ij} = \sum_{l=1}^h a_{il}c_{lj}BD = (A)_i \cdot (C)_j BD = (AC)_{ij}BD$$

The result now follows since

$$AC \otimes BD = \begin{bmatrix} (AC)_{11}BD & \cdots & (AC)_{1n}BD \\ \vdots & & \vdots \\ (AC)_{m1}BD & \cdots & (AC)_{mn}BD \end{bmatrix} \quad \square$$

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 7.8.** 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 (7.9) when  $n = m$ , we see that

$$\text{tr}(A \otimes B) = \sum_{i=1}^m a_{ii} \text{tr}(B) = \left( \sum_{i=1}^m a_{ii} \right) \text{tr}(B) = \text{tr}(A) \text{tr}(B),$$

so that the result holds. □

Theorem 7.8 gives a simplified expression for the trace of a Kronecker product. There is an analogous result for the determinant of a Kronecker product. But before we get to that, 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 7.9.** 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 7.7, 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. □

**Theorem 7.10.** 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 there exist nonsingular matrices  $P$  and  $Q$  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 since  $T_1$  and  $T_2$  are upper triangular. The result now follows since the eigenvalues of  $T_1 \otimes T_2$  are its diagonal elements, and these 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 7.11.** 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 the previous theorem

$$\begin{aligned} |A \otimes B| &= \prod_{j=1}^p \prod_{i=1}^m \lambda_i \theta_j = \prod_{j=1}^p \theta_j^m \left( \prod_{i=1}^m \lambda_i \right) = \prod_{j=1}^p \theta_j^m |A| \\ &= |A|^p \left( \prod_{j=1}^p \theta_j \right)^m = |A|^p |B|^m \end{aligned} \quad \square$$

Our final result on Kronecker products identifies a relationship between  $\text{rank}(A \otimes B)$ , and  $\text{rank}(A)$  and  $\text{rank}(B)$ .

**Theorem 7.12.** 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 utilizes Theorem 3.11, 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 2.10 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 7.10, 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  $\lambda_i$ s times the number of nonzero  $\theta_j$ s. But, since  $AA'$  and  $BB'$  are symmetric, the number of nonzero  $\lambda_i$ s is given by  $\text{rank}(AA') = \text{rank}(A)$ , and the number of nonzero  $\theta_j$ s is given by  $\text{rank}(BB') = \text{rank}(B)$ . The proof is now complete.  $\square$

**Example 7.4.** 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 Examples 3.14, 6.10, and 6.11. 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

$$y = X\beta + \epsilon,$$

where  $X = (\mathbf{1}_k \otimes \mathbf{1}_n, \mathbf{I}_k \otimes \mathbf{1}_n)$ ,  $\beta = (\mu, \tau_1, \dots, \tau_k)'$ ,  $y = (y'_1, \dots, y'_k)'$ , and  $y_i = (y_{i1}, \dots, y_{in})'$ . Consequently, a least squares solution for  $\beta$  is easily computed as

$$\begin{aligned} \hat{\beta} &= (X'X)^{-1}X'y = \left\{ \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ \mathbf{I}_k \otimes \mathbf{1}'_n \end{bmatrix} [\mathbf{1}_k \otimes \mathbf{1}_n \quad \mathbf{I}_k \otimes \mathbf{1}_n] \right\}^{-1} \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ \mathbf{I}_k \otimes \mathbf{1}'_n \end{bmatrix} y \\ &= \begin{bmatrix} nk & n\mathbf{1}'_k \\ n\mathbf{1}_k & n\mathbf{I}_k \end{bmatrix}^{-1} \begin{bmatrix} \mathbf{1}'_k \otimes \mathbf{1}'_n \\ \mathbf{I}_k \otimes \mathbf{1}'_n \end{bmatrix} y \end{aligned}$$

$$\begin{aligned}
&= \begin{bmatrix} (nk)^{-1} & \mathbf{0}' \\ \mathbf{0} & n^{-1}(\mathbf{I}_k - k^{-1}\mathbf{1}_k\mathbf{1}_k') \end{bmatrix} \begin{bmatrix} \mathbf{1}_k' \otimes \mathbf{1}_n' \\ \mathbf{I}_k \otimes \mathbf{1}_n' \end{bmatrix} \mathbf{y} \\
&= \begin{bmatrix} (nk)^{-1}(\mathbf{1}_k' \otimes \mathbf{1}_n') \\ n^{-1}(\mathbf{I}_k \otimes \mathbf{1}_n') - (nk)^{-1}(\mathbf{1}_k\mathbf{1}_k' \otimes \mathbf{1}_n') \end{bmatrix} \mathbf{y}
\end{aligned}$$

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 since  $X$  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

$$\begin{aligned}
(\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}) &= \mathbf{y}'(\mathbf{I}_{nk} - X(X'X)^-X')\mathbf{y} = \mathbf{y}'(\mathbf{I}_{nk} - n^{-1}(\mathbf{I}_k \otimes \mathbf{1}_n\mathbf{1}_n'))\mathbf{y} \\
&= \sum_{i=1}^k \mathbf{y}_i'(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n')\mathbf{y}_i = \sum_{i=1}^k \sum_{j=1}^n (y_{ij} - \bar{y}_i)^2
\end{aligned}$$

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  $\hat{\mu} = \bar{y}$ , while 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 7.5.** 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 6.17). Here  $\mu$  can be described as an overall effect, while  $\tau_i$  is an effect due to the  $i$ th level of factor A,  $\gamma_j$  is an effect due to the  $j$ th level of factor B, and  $\eta_{ij}$  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,  $\beta = (\mu, \tau_1, \dots, \tau_a, \gamma_1, \dots, \gamma_b, \eta_{11}, \eta_{12}, \dots, \eta_{ab-1}, \eta_{ab})'$  and the response vector,  $y = (y_{111}, \dots, y_{11n}, y_{121}, \dots, y_{1bn}, y_{211}, \dots, y_{abn})'$ , then the model above can be written as

$$y = X\beta + \epsilon,$$

where

$$X = (\mathbf{1}_a \otimes \mathbf{1}_b \otimes \mathbf{1}_n, \mathbf{I}_a \otimes \mathbf{1}_b \otimes \mathbf{1}_n, \mathbf{1}_a \otimes \mathbf{I}_b \otimes \mathbf{1}_n, \mathbf{I}_a \otimes \mathbf{I}_b \otimes \mathbf{1}_n).$$

Now it is easily verified that the matrix

$$X'X = \begin{bmatrix} abn & bn\mathbf{1}'_a & an\mathbf{1}'_b & n\mathbf{1}'_a \otimes \mathbf{1}'_b \\ bn\mathbf{1}_a & bn\mathbf{I}_a & n\mathbf{1}_a \otimes \mathbf{1}'_b & n\mathbf{I}_a \otimes \mathbf{1}'_b \\ an\mathbf{1}_b & n\mathbf{1}'_a \otimes \mathbf{1}_b & an\mathbf{I}_b & n\mathbf{1}'_a \otimes \mathbf{I}_b \\ n\mathbf{1}_a \otimes \mathbf{1}_b & n\mathbf{I}_a \otimes \mathbf{1}_b & n\mathbf{1}_a \otimes \mathbf{I}_b & n\mathbf{I}_a \otimes \mathbf{I}_b \end{bmatrix}$$

has as a generalized inverse the matrix

$$\text{diag}((abn)^{-1}, (bn)^{-1}(\mathbf{I}_a - a^{-1}\mathbf{1}_a\mathbf{1}'_a), (an)^{-1}(\mathbf{I}_b - b^{-1}\mathbf{1}_b\mathbf{1}'_b), C),$$

where

$$C = n^{-1}\mathbf{I}_a \otimes \mathbf{I}_b - (bn)^{-1}\mathbf{I}_a \otimes \mathbf{1}_b\mathbf{1}'_b - (an)^{-1}\mathbf{1}_a\mathbf{1}'_a \otimes \mathbf{I}_b \\ + (abn)^{-1}\mathbf{1}_a\mathbf{1}'_a \otimes \mathbf{1}_b\mathbf{1}'_b$$

Using this generalized inverse, we find that a least squares solution for  $\beta$  is given by

$$\hat{\beta} = (X'X)^{-1}X'y = \begin{bmatrix} \bar{y}_{..} \\ \bar{y}_{1.} - \bar{y}_{..} \\ \vdots \\ \bar{y}_{a.} - \bar{y}_{..} \\ \bar{y}_{.1} - \bar{y}_{..} \\ \vdots \\ \bar{y}_{.b} - \bar{y}_{..} \\ \bar{y}_{11} - \bar{y}_{1.} - \bar{y}_{.1} + \bar{y}_{..} \\ \vdots \\ \bar{y}_{ab} - \bar{y}_{a.} - \bar{y}_{.b} + \bar{y}_{..} \end{bmatrix},$$

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.

#### 4. 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

$$\text{diag}(A_1, \dots, A_r) = \begin{bmatrix} A_1 & (0) & \cdots & (0) \\ (0) & A_2 & \cdots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \cdots & A_r \end{bmatrix},$$

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



$$\text{diag}(A_1, \dots, A_r) = A_1 \oplus \dots \oplus A_r$$

Clearly, the commutative property does not hold for the direct sum since, 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  $A_1 = \dots = A_r = A$ , then

$$A_1 \oplus \dots \oplus A_r = A \oplus \dots \oplus A = I_r \otimes A$$

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 7.13.** Let  $A_1, \dots, A_r$  be matrices, where  $A_i$  is  $m_i \times m_i$ . Then

- (a)  $\text{tr}(A_1 \oplus \dots \oplus A_r) = \text{tr}(A_1) + \dots + \text{tr}(A_r)$ ,
- (b)  $|A_1 \oplus \dots \oplus A_r| = |A_1| \cdots |A_r|$ ,
- (c) if each  $A_i$  is nonsingular,  $A = A_1 \oplus \dots \oplus A_r$  is also nonsingular and  $A^{-1} = A_1^{-1} \oplus \dots \oplus A_r^{-1}$ ,
- (d)  $\text{rank}(A_1 \oplus \dots \oplus A_r) = \text{rank}(A_1) + \dots + \text{rank}(A_r)$ ,
- (e) if the eigenvalues of  $A_i$  are denoted by  $\lambda_{i,1}, \dots, \lambda_{i,m_i}$ , the eigenvalues of  $A_1 \oplus \dots \oplus A_r$  are given by  $\{\lambda_{i,j}; i = 1, \dots, r; j = 1, \dots, m_i\}$ .

## 5. 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  $a_i$  as its  $i$ th column, then  $\text{vec}(A)$  is the  $mn \times 1$  vector given by

$$\text{vec}(A) = \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_n \end{bmatrix}$$

**Example 7.6.** If  $A$  is the  $2 \times 3$  matrix given by

$$A = \begin{bmatrix} 2 & 0 & 5 \\ 8 & 1 & 3 \end{bmatrix},$$

then  $\text{vec}(A)$  is the  $6 \times 1$  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

$$\text{vec}(\mathbf{ab}') = \text{vec}([b_1\mathbf{a}, b_2\mathbf{a}, \dots, b_n\mathbf{a}]) = \begin{bmatrix} b_1\mathbf{a} \\ b_2\mathbf{a} \\ \vdots \\ b_n\mathbf{a} \end{bmatrix} = \mathbf{b} \otimes \mathbf{a}$$

Our first theorem gives this result and some others that follow directly from the definition of the  $\text{vec}$  operator.

**Theorem 7.14.** Let  $\mathbf{a}$  and  $\mathbf{b}$  be any two vectors, while  $A$  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  $\alpha$  and  $\beta$  are scalars.

The trace of a product of two matrices can be expressed in terms of the  $\text{vecs}$  of those two matrices. This result is given next.

**Theorem 7.15.** 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  $a_1, \dots, a_n$  denote the columns of  $A$  and  $b_1, \dots, b_n$  denote the columns of  $B$ . Then

$$\begin{aligned} \text{tr}(A'B) &= \sum_{i=1}^n (A'B)_{ii} = \sum_{i=1}^n a_i' b_i = [a_1', \dots, a_n'] \begin{bmatrix} b_1 \\ \vdots \\ b_n \end{bmatrix} \\ &= \{\text{vec}(A)\}' \text{vec}(B) \end{aligned} \quad \square$$

A generalization of Theorem 7.14(b) to the situation involving the  $\text{vec}$  of the product of three matrices is our next result.

**Theorem 7.16.** 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  $b_1, \dots, b_p$  are the columns of  $B$ , then  $B$  can be written as

$$B = \sum_{i=1}^p b_i e_i'$$

where  $e_i$  is the  $i$ th column of  $I_p$ . Thus,

$$\begin{aligned} \text{vec}(ABC) &= \text{vec} \left\{ A \left( \sum_{i=1}^p b_i e_i' \right) C \right\} = \sum_{i=1}^p \text{vec}(A b_i e_i' C) \\ &= \sum_{i=1}^p \text{vec} \{ (A b_i) (C' e_i)' \} = \sum_{i=1}^p C' e_i \otimes A b_i \\ &= (C' \otimes A) \sum_{i=1}^p (e_i \otimes b_i), \end{aligned}$$

where the second last equality follows from Theorem 7.14(b). The result now follows since, by again using Theorem 7.14(b), we find that

$$\sum_{i=1}^p (\mathbf{e}_i \otimes \mathbf{b}_i) = \sum_{i=1}^p \text{vec}(\mathbf{b}_i \mathbf{e}_i') = \text{vec} \left( \sum_{i=1}^p \mathbf{b}_i \mathbf{e}_i' \right) = \text{vec}(B) \quad \square$$

**Example 7.7.** In Chapter 6, we discussed systems of linear equations of the form  $A\mathbf{x} = \mathbf{c}$ , as well as systems of equations of the form  $AXB = C$ . Using the  $\text{vec}$  operator and Theorem 7.16, this second system of equations can be equivalently expressed as

$$\text{vec}(AXB) = (B' \otimes A)\text{vec}(X) = \text{vec}(C);$$

that is, this is an equation of the form  $A\mathbf{x} = \mathbf{c}$ , where in place of  $A$ ,  $\mathbf{x}$ , and  $\mathbf{c}$ , we have  $(B' \otimes A)$ ,  $\text{vec}(X)$ , and  $\text{vec}(C)$ . As a result, Theorem 6.4, which gives the general form of a solution to  $A\mathbf{x} = \mathbf{c}$ , can be used to prove Theorem 6.5, which gives the general form of a solution to  $AXB = C$ . The details of this proof are left to the reader.

Theorem 7.15 also can be generalized to a result involving the product of more than two matrices.

**Theorem 7.17.** 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 7.15, it follows that

$$\text{tr}(ABCD) = \text{tr}\{A(BCD)\} = \{\text{vec}(A')\}' \text{vec}(BCD)$$

But from the previous theorem, 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 7.17 are left to the reader as an exercise.

**Corollary 7.17.1.** Let  $A$  and  $C$  be matrices of sizes  $m \times n$  and  $n \times m$ , respectively, while  $B$  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 will be 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  $\tilde{v}(A)$  since the diagonal elements of  $A$  must be zero, while  $a_{ji} = -a_{ij}$  if  $i \neq j$ . The notation we use here, that is,  $v(A)$  and  $\tilde{v}(A)$ , corresponds to that used by Magnus (1988). Others [see, for example, Henderson and Searle (1979)] use the notation  $\text{vech}(A)$  and  $\text{veck}(A)$ . In Section 8, we will discuss some transformations which relate the  $v$  and  $\tilde{v}$  operators to the  $\text{vec}$  operator.

**Example 7.8.** 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 results 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 1s from  $v(R)$ , we obtain

$$\tilde{v}(R) = (r_{12}, r_{13}, r_{23})',$$

which contains all of the random variables in  $R$ .

## 6. THE HADAMARD PRODUCT

A matrix operator that is a little more obscure than our other matrix operators, but one which 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 elementwise 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 7.9.** If  $A$  and  $B$  are the  $2 \times 3$  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 9.7. 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, the reader is referred to Styan (1973) and Horn and Johnson (1991). We begin with some elementary properties that follow directly from the definition of the Hadamard product.

**Theorem 7.18.** 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)  $C(A \odot B) = (CA) \odot B = A \odot (CB)$  and  $(A \odot B)C = (AC) \odot B = A \odot (BC)$ ,  
 if  $m = n$  and  $C$  is diagonal,  
 (i)  $\mathbf{ab}' \odot \mathbf{cd}' = (\mathbf{a} \odot \mathbf{c})(\mathbf{b} \odot \mathbf{d})'$ , 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  $\Psi_m$  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  $A \otimes B$ , containing 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 7.19.** Let  $A$  and  $B$  be  $m \times n$  matrices. Then

$$\text{rank}(A \odot B) \leq \text{rank}(A) \text{rank}(B)$$

*Proof.* Let  $r_A = \text{rank}(A)$  and  $r_B = \text{rank}(B)$ . It follows from the singular value decomposition theorem (Theorem 4.1 and Corollary 4.1.1) that there exist  $m \times r_A$  and  $n \times r_A$  matrices  $U = (\mathbf{u}_1, \dots, \mathbf{u}_{r_A})$  and  $V = (\mathbf{v}_1, \dots, \mathbf{v}_{r_A})$ , and  $m \times r_B$  and  $n \times r_B$  matrices  $W = (\mathbf{w}_1, \dots, \mathbf{w}_{r_B})$  and  $X = (\mathbf{x}_1, \dots, \mathbf{x}_{r_B})$ , such that  $A = UV'$  and  $B = WX'$ . Then

$$\begin{aligned}
 A \odot B &= UV' \odot WX' = \left( \sum_{i=1}^{r_A} u_i v_i' \right) \odot \left( \sum_{j=1}^{r_B} w_j x_j' \right) \\
 &= \sum_{i=1}^{r_A} \sum_{j=1}^{r_B} (u_i v_i' \odot w_j x_j') = \sum_{i=1}^{r_A} \sum_{j=1}^{r_B} (u_i \odot w_j)(v_i \odot x_j)',
 \end{aligned}$$

where we have used Theorem 7.18(c) and (i). The result follows since we have now expressed  $A \odot B$  as the sum of  $r_A r_B$  matrices, each having rank of at most one.  $\square$

**Example 7.10.** While Theorem 7.19 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 while  $A \odot B$  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, yet  $A \odot B$  has rank 0 since  $A \odot B = (0)$ .

The following result shows that a bilinear form in a Hadamard product of two matrices may be written as a trace.

**Theorem 7.20.** Let  $A$  and  $B$  be  $m \times n$  matrices, and let  $x$  and  $y$  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)  $x'(A \odot B)y = \text{tr}(D_x A D_y B')$ ,

where  $D_x = \text{diag}(x_1, \dots, x_m)$  and similarly for  $D_y$ .

*Proof.* (a) follows since

$$\begin{aligned}
 \mathbf{1}_m'(A \odot B)\mathbf{1}_n &= \sum_{i=1}^m \sum_{j=1}^n (A \odot B)_{ij} = \sum_{i=1}^m \sum_{j=1}^n a_{ij} b_{ij} \\
 &= \sum_{i=1}^m (A)_{i \cdot} (B')_{\cdot i} = \sum_{i=1}^m (AB')_{ii} = \text{tr}(AB')
 \end{aligned}$$



To prove (b), note that  $x = D_x \mathbf{1}_m$  and  $y = D_y \mathbf{1}_n$ , so that by using Theorem 7.20(a) and Theorem 7.18(h), we find that

$$x'(A \odot B)y = \mathbf{1}'_m D_x (A \odot B) D_y \mathbf{1}_n = \mathbf{1}'_m (D_x A \odot B D_y) \mathbf{1}_n = \text{tr}(D_x A D_y B') \quad \square$$

The following result can be helpful in determining whether the Hadamard product of two symmetric matrices is nonnegative definite or positive definite.

**Theorem 7.21.** 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 also 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  $k$  since  $B$  is nonnegative definite. Then we find that for any  $m \times 1$  vector  $y$ ,

$$\begin{aligned} y'(A \odot B)y &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ij} y_i y_j = \sum_{k=1}^m \left( \sum_{i=1}^m \sum_{j=1}^m \lambda_k (y_i x_{ik}) a_{ij} (y_j x_{jk}) \right) \\ &= \sum_{k=1}^m \lambda_k (y \odot x_k)' A (y \odot x_k), \end{aligned} \quad (7.11)$$

where  $x_k$  represents the  $k$ th column of  $X$ . Since  $A$  is nonnegative definite, the sum in (7.11) must be nonnegative, and so  $A \odot B$  is also nonnegative definite. This proves (a). Now if  $A$  is positive definite, then (7.11) will be positive for any  $y \neq \mathbf{0}$  that satisfies  $y \odot x_k \neq \mathbf{0}$  for at least one  $k$  for which  $\lambda_k > 0$ . But if  $B$  is also positive definite, then  $\lambda_k > 0$  for all  $k$  and if  $y$  has its  $h$ th component  $y_h \neq 0$ , then  $y \odot x_k = \mathbf{0}$  for all  $k$  only if the  $h$ th row of  $X$  has all zeros. This is not possible since  $X$  is nonsingular. Consequently, there is no  $y \neq \mathbf{0}$  for which (7.11) equals zero and so (b) follows.  $\square$

Theorem 7.21(b) gives a sufficient condition for the matrix  $A \odot B$  to be positive definite. The following example demonstrates that this condition is not necessary.

**Example 7.11.** Consider the  $2 \times 2$  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 since, 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 since  $A \odot B = B$ . However,  $A$  is not positive definite since  $\text{rank}(A) = 1$ .

A sufficient condition for the positive definiteness of  $A \odot B$ , weaker than that given in Theorem 7.21(b), is given in our next theorem.

**Theorem 7.22.** 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  $x \neq 0$ ,  $x'(A \odot B)x > 0$ . Since  $B$  is positive definite, there exists a nonsingular matrix  $T$  such that  $B = TT'$ . It follows then from Theorem 7.20(b) that

$$x'(A \odot B)x = \text{tr}(D_x A D_x B') = \text{tr}(D_x A D_x T T') = \text{tr}(T' D_x A D_x T) \quad (7.12)$$

Since  $A$  is nonnegative definite, so is  $D_x A D_x$ . In addition, if  $x \neq 0$ , and  $A$  has no diagonal elements equal to zero, then  $D_x A D_x \neq (0)$ ; that is,  $D_x A D_x$  has rank of at least one, and so it has at least one positive eigenvalue. Since  $T$  is nonsingular,  $\text{rank}(D_x A D_x) = \text{rank}(T' D_x A D_x T)$ , and so  $T' D_x A D_x T$  is also nonnegative definite with at least one positive eigenvalue. The result now follows since (7.12) implies that  $x'(A \odot B)x$  is the sum of the eigenvalues of  $T' D_x A D_x T$ .  $\square$

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 7.23.** 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

$$\begin{aligned}
 |A| &= a_{11} \begin{vmatrix} a_{22} & a_{23} & \cdots & a_{2m} \\ a_{32} & a_{33} & \cdots & a_{3m} \\ \vdots & \vdots & & \vdots \\ a_{m2} & a_{m3} & \cdots & a_{mm} \end{vmatrix} + \begin{vmatrix} 0 & a_{12} & \cdots & a_{1m} \\ a_{21} & a_{22} & \cdots & a_{2m} \\ \vdots & \vdots & & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mm} \end{vmatrix} \\
 &= a_{11}|A_1| + \begin{vmatrix} 0 & \mathbf{a}' \\ \mathbf{a} & A_1 \end{vmatrix}, \tag{7.13}
 \end{aligned}$$

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 (7.13), 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 7.23.1.** 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

$$a_{ii} = (BB')_{ii} = (B)_i \cdot (B')_{\cdot i} = (B)_i \cdot (B)_i' = \sum_{j=1}^m b_{ij}^2,$$

and so the result follows immediately from Theorem 7.23.  $\square$

Theorem 7.23 also holds for positive semidefinite matrices except that in this case  $A$  need not be diagonal for equality since one or more of its diagonal elements may equal zero. Likewise, Corollary 7.23.1 holds for singular matrices except for the statement concerning equality.

Hadamard's inequality given in Theorem 7.23 can be expressed, using the Hadamard product, as

$$|A| \left( \prod_{i=1}^m 1 \right) \leq |A \odot I_m|, \quad (7.14)$$

where the term  $(\prod 1)$  corresponds to the product of the diagonal elements of  $I_m$ . Theorem 7.25 will show that the inequality (7.14) holds for other matrices besides the identity. But first we will need the following result.

**Theorem 7.24.** 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  $A_\alpha$  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 since  $A$  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,  $A_\alpha$  may be written as

$$\begin{aligned}
 A_\alpha &= \begin{bmatrix} a_{11} & a' \\ \mathbf{a} & A_1 \end{bmatrix} - \begin{bmatrix} (a_{11} - a'A_1^{-1}a) & \mathbf{0}' \\ \mathbf{0} & (0) \end{bmatrix} \\
 &= \begin{bmatrix} a'A_1^{-1}a & a' \\ \mathbf{a} & A_1 \end{bmatrix} = \begin{bmatrix} a'A_1^{-1} \\ I_{m-1} \end{bmatrix} A_1 [A_1^{-1}a \quad I_{m-1}]
 \end{aligned}$$

Since  $A_1$  is positive definite, there exists an  $(m - 1) \times (m - 1)$  matrix  $T$  such that  $A_1 = TT'$ . If we let  $U' = T'[A_1^{-1}a \quad I_{m-1}]$ , then  $A_\alpha = UU'$ , and so  $A_\alpha$  is nonnegative definite.  $\square$

**Theorem 7.25.** 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 since  $|A| = 0$ , while  $|A \odot B| \geq 0$  is guaranteed by Theorem 7.21. For the case in which  $A$  is positive definite, we will prove the result by induction. The result holds when  $m = 2$ , since 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

$$|A_1| \prod_{i=2}^m b_{ii} \leq |A_1 \odot B_1|, \tag{7.15}$$

where  $A_1$  and  $B_1$  are the submatrices of  $A$  and  $B$  formed by deleting their first row and first column. From Theorem 7.24 we know that  $(A - \alpha e_1 e_1')$  is nonnegative definite, where  $\alpha = |A|/|A_1|$ . Thus, by using Theorem 7.21(a), Theorem 7.18(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. But  $(A \odot B)_1 = A_1 \odot B_1$  so that the inequality above, along with (7.15) 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 7.26.** Let  $A$  and  $B$  be  $m \times m$  symmetric matrices. If  $B$  is nonnegative 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 there exists an  $m \times m$  matrix  $T$  such that  $B = TT'$ . Let  $t_j$  be the  $j$ th column of  $T$ , while  $t_{ij}$  denotes the  $(i, j)$ th element of  $T$ . For any  $m \times 1$  vector,  $x \neq 0$ , we find that

$$\begin{aligned} x'(A \odot B)x &= \sum_{i=1}^m \sum_{j=1}^m a_{ij} b_{ij} x_i x_j = \sum_{i=1}^m \sum_{j=1}^m a_{ij} \left( \sum_{h=1}^m t_{ih} t_{jh} \right) x_i x_j \\ &= \sum_{h=1}^m \left( \sum_{i=1}^m \sum_{j=1}^m (x_i t_{ih}) a_{ij} (x_j t_{jh}) \right) = \sum_{h=1}^m (x \odot t_h)' A (x \odot t_h) \\ &\leq \lambda_1(A) \sum_{h=1}^m (x \odot t_h)' (x \odot t_h) = \lambda_1(A) \sum_{h=1}^m \sum_{j=1}^m x_j^2 t_{jh}^2 \\ &= \lambda_1(A) \sum_{j=1}^m x_j^2 \left( \sum_{h=1}^m t_{jh}^2 \right) = \lambda_1(A) \sum_{j=1}^m x_j^2 b_{jj} \\ &\leq \lambda_1(A) \left\{ \max_{1 \leq i \leq m} b_{ii} \right\} x'x, \end{aligned} \tag{7.16}$$

where the first inequality arises from the relation

$$\lambda_1(A) = \max_{y'y \neq 0} \frac{y'Ay}{y'y}$$

given in Theorem 3.15. Using this same relationship for  $A \odot B$ , along with (7.16), we find that for any  $i$ ,  $1 \leq i \leq m$ ,

$$\lambda_i(A \odot B) \leq \lambda_1(A \odot B) = \max_{x'x \neq 0} \frac{x'(A \odot B)x}{x'x} \leq \lambda_1(A) \left\{ \max_{1 \leq i \leq m} b_{ii} \right\},$$

which is the required upper bound on  $\lambda_i(A \odot B)$ . The lower bound is obtained in a similar fashion by using the identity

$$\lambda_m(A) = \min_{y'y \neq 0} \frac{y'Ay}{y'y} \quad \square$$

Our final result provides an alternative lower bound for the eigenvalues of  $(A \odot B)$ . The derivation of this bound will make use of the following result.

**Theorem 7.27.** Let  $A$  be an  $m \times m$  positive definite matrix. Then the matrix  $(A \odot A^{-1}) - I_m$  is nonnegative definite.

*Proof.* Let  $\sum_{i=1}^m \lambda_i x_i x_i'$  be the spectral decomposition of  $A$  so that  $A^{-1} = \sum_{i=1}^m \lambda_i^{-1} x_i x_i'$ . Then

$$\begin{aligned} (A \odot A^{-1}) - I_m &= (A \odot A^{-1}) - I_m \odot I_m \\ &= \left( \sum_{i=1}^m \lambda_i x_i x_i' \odot \sum_{j=1}^m \lambda_j^{-1} x_j x_j' \right) - \left( \sum_{i=1}^m x_i x_i' \odot \sum_{j=1}^m x_j x_j' \right) \\ &= \sum_{i=1}^m \sum_{j=1}^m (\lambda_i \lambda_j^{-1} - 1) (x_i x_i' \odot x_j x_j') \\ &= \sum_{i \neq j} (\lambda_i \lambda_j^{-1} - 1) (x_i x_i' \odot x_j x_j') \\ &= \sum_{i < j} (\lambda_i \lambda_j^{-1} + \lambda_j \lambda_i^{-1} - 2) (x_i \odot x_j) (x_i \odot x_j)' = XDX', \end{aligned}$$

where  $X$  is the  $m \times m(m-1)/2$  matrix having  $(x_i \odot x_j)$ ,  $i < j$  as its columns, while  $D$  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  $XD X'$ .  $\square$

**Theorem 7.28.** Let  $A$  and  $B$  be  $m \times m$  nonnegative definite matrices. Then

$$\lambda_m(A \odot B) \geq \lambda_m(AB)$$

*Proof.* Due to Theorem 7.21,  $A \odot B$  is nonnegative definite, and so the inequality is obvious if either  $A$  or  $B$  is singular since 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  $TAT' - \lambda_m(AB)I_m$  is nonnegative definite since its  $i$ th largest eigenvalue is  $\lambda_i(TAT') - \lambda_m(AB)$ , and  $\lambda_m(AB) = \lambda_m(TAT')$ . 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 due to Theorem 7.21, while  $\lambda_m(AB)\{(B^{-1} \odot B) - I_m\}$  is nonnegative definite due to Theorem 7.27, 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  $x$ ,

$$x'(A \odot B)x \geq \lambda_m(AB)x'x,$$

and so the result follows from Theorem 3.15.  $\square$

## 7. 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 7.1.** 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

$$K_{mn} = \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij}) \tag{7.17}$$

The matrix  $H_{ij}$  can be conveniently expressed in terms of columns from the identity matrices  $I_m$  and  $I_n$ . If  $e_{i,m}$  is the  $i$ th column of  $I_m$  and  $e_{j,n}$  is the  $j$ th column of  $I_n$ , then  $H_{ij} = e_{i,m}e'_{j,n}$ .

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 (7.17), it is easy to verify that  $K_{16} = K_{61} = I_6$ , while

$$K_{23} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix},$$

$$K_{32} = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

The fact that  $K_{32} = K'_{23}$  is not a coincidence, since this is a general property that follows from the definition of  $K_{mn}$ .

**Theorem 7.29.** The commutation matrix satisfies the properties

- (a)  $K_{m1} = K_{1m} = I_m$ ,
- (b)  $K'_{mn} = K_{nm}$ ,
- (c)  $K^{-1}_{mn} = K_{nm}$ .

*Proof.* When  $H_{ij}$  is  $m \times 1$ , then  $H_{ij} = e_{i,m}$  and so

$$K_{m1} = \sum_{i=1}^m (e_{i,m} \otimes e'_{i,m}) = I_m = \sum_{i=1}^m (e'_{i,m} \otimes e_{i,m}) = K_{1m},$$

proving (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 since

$$H_{ij}H'_{kl} = e_{i,m}e'_{j,n}e_{l,n}e'_{k,m} = \begin{cases} e_{i,m}e'_{k,m}, & \text{if } j = l, \\ (0), & \text{if } j \neq l, \end{cases}$$

$$H'_{ij}H_{kl} = e_{j,n}e'_{i,m}e_{k,m}e'_{l,n} = \begin{cases} e_{j,n}e'_{l,n}, & \text{if } i = k, \\ (0), & \text{if } i \neq k, \end{cases}$$

and so

$$\begin{aligned} K_{mn}K_{nm} &= K_{mn}K'_{mn} = \left\{ \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij}) \right\} \left\{ \sum_{k=1}^m \sum_{l=1}^n (H_{kl} \otimes H'_{kl})' \right\} \\ &= \sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^m \sum_{l=1}^n (H_{ij}H'_{kl} \otimes H'_{ij}H_{kl}) \\ &= \sum_{i=1}^m \sum_{j=1}^n (e_{i,m}e'_{i,m} \otimes e_{j,n}e'_{j,n}) \\ &= \left( \sum_{i=1}^m e_{i,m}e'_{i,m} \right) \otimes \left( \sum_{j=1}^n e_{j,n}e'_{j,n} \right) = I_m \otimes I_n = I_{mn} \quad \square \end{aligned}$$

Commutation matrices have important relationships with the  $\text{vec}$  operator and the Kronecker product. For an  $m \times n$  matrix  $A$ , the two vectors  $\text{vec}(A)$  and  $\text{vec}(A')$  are related since 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 which transforms  $\text{vec}(A)$  to  $\text{vec}(A')$ .

**Theorem 7.30.** For any  $m \times n$  matrix  $A$ ,

$$K_{mn} \text{vec}(A) = \text{vec}(A')$$

*Proof.* Clearly, since  $a_{ij}H'_{ij}$  is the  $n \times m$  matrix whose only nonzero element,  $a_{ij}$ , is in the  $(j, i)$ th position, we have

$$\begin{aligned} A' &= \sum_{i=1}^m \sum_{j=1}^n a_{ij}H'_{ij} = \sum_{i=1}^m \sum_{j=1}^n (e'_{i,m}Ae_{j,n})e_{j,n}e'_{i,m} \\ &= \sum_{i=1}^m \sum_{j=1}^n e_{j,n}(e'_{i,m}Ae_{j,n})e'_{i,m} = \sum_{i=1}^m \sum_{j=1}^n (e_{j,n}e'_{i,m})A(e_{j,n}e'_{i,m}) \\ &= \sum_{i=1}^m \sum_{j=1}^n H'_{ij}AH'_{ij} \end{aligned}$$

The result now follows by taking the  $\text{vec}$  of both sides and using Theorem 7.16, since

$$\begin{aligned} \text{vec}(A') &= \text{vec}\left(\sum_{i=1}^m \sum_{j=1}^n H'_{ij}AH'_{ij}\right) = \sum_{i=1}^m \sum_{j=1}^n \text{vec}(H'_{ij}AH'_{ij}) \\ &= \sum_{i=1}^m \sum_{j=1}^n (H_{ij} \otimes H'_{ij})\text{vec}(A) = K_{mn}\text{vec}(A) \quad \square \end{aligned}$$

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 7.31.

**Theorem 7.31.** Let  $A$  be an  $m \times n$  matrix,  $B$  be a  $p \times q$  matrix,  $x$  be an  $m \times 1$  vector, and  $y$  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)  $K_{pm}(A \otimes y) = y \otimes A$ ,
- (d)  $K_{mp}(y \otimes A) = A \otimes y$ ,
- (e)  $K_{pm}(x \otimes y) = y \otimes x$ ,
- (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 Theorems 7.16 and 7.30, 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 using Theorem 7.29(c) yields (b). Properties (c)–(e) follow from (a) and Theorem 7.29(a) since

$$\begin{aligned} K_{pm}(A \otimes y) &= (y \otimes A)K_{1n} = y \otimes A, \\ K_{mp}(y \otimes A) &= (A \otimes y)K_{n1} = A \otimes y, \\ K_{pm}(x \otimes y) &= (y \otimes x)K_{11} = y \otimes x \end{aligned}$$

To prove (f), use the definition of the commutation matrix to get

$$\begin{aligned} \text{tr}\{(B \otimes A)K_{mn}\} &= \sum_{i=1}^m \sum_{j=1}^n \text{tr}\{(B \otimes A)(H_{ij} \otimes H'_{ij})\} \\ &= \sum_{i=1}^m \sum_{j=1}^n \{\text{tr}(BH_{ij})\}\{\text{tr}(AH'_{ij})\} \\ &= \sum_{i=1}^m \sum_{j=1}^n (e'_{j,n} B e_{i,m})(e'_{i,m} A e_{j,n}) = \sum_{i=1}^m \sum_{j=1}^n b_{ji} a_{ij} \\ &= \sum_{j=1}^n (B)_{j \cdot} (A)_{\cdot j} = \sum_{j=1}^n (BA)_{jj} = \text{tr}(BA) \quad \square \end{aligned}$$

The commutation matrix also can be utilized to obtain a relationship between the  $\text{vec}$  of a Kronecker product and the Kronecker product of the corresponding  $\text{vecs}$ .

**Theorem 7.32.** 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  $a_1, \dots, a_n$  be the columns of  $A$  and  $b_1, \dots, b_q$  the columns of  $B$ . Then since  $A$  and  $B$  can be written as

$$A = \sum_{i=1}^n a_i e'_{i,n}, \quad B = \sum_{j=1}^q b_j e'_{j,q},$$

we have

$$\begin{aligned} \text{vec}(A \otimes B) &= \sum_{i=1}^n \sum_{j=1}^q \text{vec}(a_i e'_{i,n} \otimes b_j e'_{j,q}) \\ &= \sum_{i=1}^n \sum_{j=1}^q \text{vec}\{(a_i \otimes b_j)(e'_{i,n} \otimes e'_{j,q})\} \\ &= \sum_{i=1}^n \sum_{j=1}^q \{(e_{i,n} \otimes e_{j,q}) \otimes (a_i \otimes b_j)\} \\ &= \sum_{i=1}^n \sum_{j=1}^q \{e_{i,n} \otimes K_{qm}(a_i \otimes e_{j,q}) \otimes b_j\} \\ &= \sum_{i=1}^n \sum_{j=1}^q (I_n \otimes K_{qm} \otimes I_p)(e_{i,n} \otimes a_i \otimes e_{j,q} \otimes b_j) \\ &= (I_n \otimes K_{qm} \otimes I_p) \left\{ \sum_{i=1}^n (e_{i,n} \otimes a_i) \right\} \otimes \left\{ \sum_{j=1}^q (e_{j,q} \otimes b_j) \right\} \\ &= (I_n \otimes K_{qm} \otimes I_p) \left\{ \sum_{i=1}^n \text{vec}(a_i e'_{i,n}) \right\} \otimes \left\{ \sum_{j=1}^q \text{vec}(b_j e'_{j,q}) \right\} \\ &= (I_n \otimes K_{qm} \otimes I_p) \{\text{vec}(A) \otimes \text{vec}(B)\} \quad \square \end{aligned}$$

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 7.33.** 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 3.8 that its eigenvalues are also real. Further, since  $K_{mm}$  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, implying that  $m^2 - p$  is the number of eigenvalues equal to +1. Since the trace equals the sum of the eigenvalues, we must have  $\text{tr}(K_{mm}) = p(-1) + (m^2 - p)(1) = m^2 - 2p$ . But by using basic properties of the trace, we also find that

$$\begin{aligned} \text{tr}(K_{mm}) &= \text{tr} \left\{ \sum_{i=1}^m \sum_{j=1}^m (e_i e_j' \otimes e_j e_i') \right\} = \sum_{i=1}^m \sum_{j=1}^m \text{tr}(e_i e_j' \otimes e_j e_i') \\ &= \sum_{i=1}^m \sum_{j=1}^m \{ \text{tr}(e_i e_j') \} \{ \text{tr}(e_j e_i') \} = \sum_{i=1}^m \sum_{j=1}^m (e_i' e_j)^2 \\ &= \sum_{i=1}^m 1 = m \end{aligned}$$

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.  $\square$

We will see later that the commutation matrix  $K_{mm}$  appears in some important matrix moment formulas through the term  $N_m = \frac{1}{2}(\mathbf{I}_{m^2} + K_{mm})$ . Consequently, we will establish some basic properties of  $N_m$ .

**Theorem 7.34.** Let  $N_m = \frac{1}{2}(\mathbf{I}_{m^2} + K_{mm})$ , 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  $\mathbf{I}_{m^2}$  and  $K_{mm}$ , while

$$N_m^2 = \frac{1}{4} (\mathbf{I}_{m^2} + K_{mm})^2 = \frac{1}{4} (\mathbf{I}_{m^2} + 2K_{mm} + K_{mm}^2) = \frac{1}{2} (\mathbf{I}_{m^2} + K_{mm}) = N_m,$$

since  $K_{mm}^2 = \mathbf{I}$  follows from the fact that  $K_{mm}^{-1} = K_{mm}$ . Similarly, (b) follows from the fact that  $K_{mm}^2 = \mathbf{I}_{m^2}$ . 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, to prove (d), note that, by using Theorem 7.31(b), and Theorem 7.34(b),

$$N_m(A \otimes B)N_m = N_m K_{mm}(B \otimes A)K_{mm}N_m = N_m(B \otimes A)N_m \quad \square$$

The proof of our final result will be left to the reader as an exercise.

**Theorem 7.35.** 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. 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  $\text{vec}$  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 since  $a_{ij} = a_{ji}$  for  $i \neq j$ . For this reason, we previously defined  $v(A)$  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  $v(A)$  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  $A$ ,

$$D_m v(A) = \text{vec}(A) \quad (7.18)$$

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 for the  $m^2 \times m(m+1)/2$  duplication matrix  $D_m$ , refer to Magnus (1988) or Problem 7.54.

Some properties of the duplication matrix and its Moore–Penrose inverse are summarized in Theorem 7.36.

**Theorem 7.36.** Let  $D_m$  be the  $m^2 \times m(m+1)/2$  duplication matrix and  $D_m^+$  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)  $D_m^+ \text{vec}(A) = v(A)$  for every  $m \times m$  symmetric matrix  $A$ .

*Proof.* Clearly, for every  $m(m+1)/2 \times 1$  vector  $x$  there exists an  $m \times m$  symmetric matrix  $A$  such that  $x = v(A)$ . But if for some symmetric  $A$ ,  $D_m v(A) = 0$ , then from the definition of  $D_m$ ,  $\text{vec}(A) = 0$ , which then implies that  $v(A) = 0$ . Thus,  $D_m x = 0$  only if  $x = 0$ , and so  $D_m$  has full column rank. Parts (b) and (c) follow immediately from (a) and Theorem 5.3, while (d) is obtained by premultiplying (7.18) by  $D_m^+$  and then using (c).  $\square$

The duplication matrix and its Moore–Penrose inverse have some important relationships with  $K_{mm}$  and  $N_m$ .

**Theorem 7.37.** Let  $D_m$  be the  $m^2 \times m(m+1)/2$  duplication matrix and  $D_m^+$  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

$$K_{mm} D_m v(A) = K_{mm} \text{vec}(A) = \text{vec}(A') = \text{vec}(A) = D_m v(A) \quad (7.19)$$

Similarly, we have

$$N_m D_m v(A) = N_m \text{vec}(A) = \frac{1}{2} \text{vec}(A + A') = \text{vec}(A) = D_m v(A) \quad (7.20)$$

Since  $\{v(A): A \text{ } m \times m \text{ and } A' = A\}$  is all of  $m(m+1)/2$ -dimensional space, (7.19) and (7.20) establish (a). To prove (b), take the transpose of (a), premultiply all three sides by  $(D_m' D_m)^{-1}$ , and then use Theorem 7.36(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

$$\begin{aligned} N_m \text{vec}(A) &= \frac{1}{2} (I_{m^2} + K_{mmm}) \text{vec}(A) = \frac{1}{2} \{ \text{vec}(A) + \text{vec}(A') \} \\ &= \text{vec}(A_*) \end{aligned}$$

Using this and (b), we find that

$$\begin{aligned} D_m D_m^+ \text{vec}(A) &= D_m D_m^+ N_m \text{vec}(A) = D_m D_m^+ \text{vec}(A_*) \\ &= D_m \text{vec}(A_*) = \text{vec}(A_*) = N_m \text{vec}(A), \end{aligned}$$

and so the proof is complete. □

We will need the following result in the next chapter

**Theorem 7.38.** 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 above yields  $I_{m(m+1)/2}$ . Using Theorem 7.35(a), Theorem 7.36(c), and Theorem 7.37(a) and (c), we have

$$\begin{aligned} &D'_m(A \otimes A)D_m D_m^+(A^{-1} \otimes A^{-1})D_m^{+'} \\ &= D'_m(A \otimes A)N_m(A^{-1} \otimes A^{-1})D_m^{+'} \\ &= D'_m N_m(A \otimes A)(A^{-1} \otimes A^{-1})D_m^{+'} = D'_m N_m D_m^{+'} \\ &= (N_m D_m)' D_m^{+'} = D'_m D_m^{+'} = (D_m^+ D_m)' = I_{m(m+1)/2} \quad \square \end{aligned}$$

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  $v(A)$  except that  $\text{vec}(A)$  has some additional zeros. We will denote by  $L'_m$  the  $m^2 \times m(m+1)/2$  matrix which transforms  $v(A)$  into  $\text{vec}(A)$ ; that is,  $L'_m$  satisfies

$$L'_m v(A) = \text{vec}(A) \tag{7.21}$$

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 following properties of the matrix  $L_m$  can be proven directly from its definition given in (7.21).

**Theorem 7.39.** 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)  $L_m \text{vec}(A) = \mathbf{v}(A)$ , for every  $m \times m$  matrix  $A$ .

*Proof.* Note that if  $A$  is lower triangular, then  $\text{vec}(A)' \text{vec}(A) = \mathbf{v}(A)' \mathbf{v}(A)$  and so (7.21) implies

$$\mathbf{v}(A)' L_m L'_m \mathbf{v}(A) - \mathbf{v}(A)' \mathbf{v}(A) = \mathbf{v}(A)' (L_m L'_m - I_{m(m+1)/2}) \mathbf{v}(A) = 0$$

for all lower triangular matrices  $A$ . But this can be true only if (b) holds since  $\{\mathbf{v}(A): A \text{ } m \times m \text{ and lower triangular}\} = R^{m(m+1)/2}$ . Part (a) follows immediately from (b), as does (c) since  $L_m^+ = (L'_m L_m)^{-1} L'_m$ . 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,

$$0 = \text{vec}(A_L)' \text{vec}(A_U) = \mathbf{v}(A_L)' L_m \text{vec}(A_U),$$

and since, 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 (7.21), (b), and the fact that  $v(A_L) = v(A)$ , we have

$$L_m \text{vec}(A) = L_m \text{vec}(A_L + A_U) = L_m \text{vec}(A_L) = L_m L'_m v(A_L) = v(A_L) = v(A) \quad \square$$

We see from property (d) in Theorem 7.39 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. Our next result 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 7.40.** 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 very similar to  $L_m$ , some of its basic properties parallel those of  $L_m$ . For instance, the following results are analogous to those in Theorem 7.39. The proofs, which we omit, are similar to those of Theorem 7.39.

**Theorem 7.41.** 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  $A$ .

Our final theorem 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 7.42.** 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)  $\tilde{L}_m L'_m L_m \tilde{L}'_m = I_{m(m-1)/2}$ .

## 9. 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. This 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 (1973), as well as the books by Gantmacher (1959) and Horn and Johnson (1985). Most of the proofs that we present here follow along the lines of the derivations, based on matrix norms, given in Horn and Johnson (1985).

We begin with some results regarding the spectral radius of nonnegative and positive matrices.

**Theorem 7.43.** Let  $A$  be an  $m \times m$  matrix and  $x$  be an  $m \times 1$  vector. If  $A \geq (0)$  and  $x > 0$ , then

$$\min_{1 \leq i \leq m} \sum_{j=1}^m a_{ij} \leq \rho(A) \leq \max_{1 \leq i \leq m} \sum_{j=1}^m a_{ij}, \quad (7.22)$$

$$\min_{1 \leq i \leq m} x_i^{-1} \sum_{j=1}^m a_{ij} x_j \leq \rho(A) \leq \max_{1 \leq i \leq m} x_i^{-1} \sum_{j=1}^m a_{ij} x_j \quad (7.23)$$

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  $k$ ,  $A^k \geq B^k$  and this then 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 4.24 that  $\rho(A) \geq \rho(B)$ . But this proves the lower bound in (7.22) since  $\rho(B) = \alpha$  follows from the fact that  $\rho(B) \geq \alpha$  since

$$B\mathbf{1}_m = \alpha\mathbf{1}_m$$

and  $\rho(B) \leq \|B\|_{\infty} = \alpha$  due to Theorem 4.19. 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 (7.23) follow directly from those in (7.22) since 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 7.44.** Let  $A$  be an  $m \times m$  positive matrix. Then  $\rho(A)$  is positive and is an eigenvalue of  $A$ . In addition, there exists a positive eigenvector of  $A$  corresponding to the eigenvalue  $\rho(A)$ .

*Proof.*  $\rho(A) > 0$  follows immediately from Theorem 7.43 since  $A$  is positive. By the definition of  $\rho(A)$ , there exists an eigenvalue of  $A$ ,  $\lambda$ , such that

$|\lambda| = \rho(A)$ . Let  $\mathbf{x}$  be an eigenvector of  $A$  corresponding to  $\lambda$  so that  $A\mathbf{x} = \lambda\mathbf{x}$ . Note that

$$\rho(A)\text{abs}(\mathbf{x}) = |\lambda| \text{abs}(\mathbf{x}) = \text{abs}(\lambda\mathbf{x}) = \text{abs}(A\mathbf{x}) \leq \text{abs}(A)\text{abs}(\mathbf{x}) = A \text{abs}(\mathbf{x}),$$

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  $\mathbf{y} = A \text{abs}(\mathbf{x}) - \rho(A)\text{abs}(\mathbf{x})$  is nonnegative. The vector  $\mathbf{z} = A \text{abs}(\mathbf{x})$  is positive since  $A$  is positive and the eigenvector  $\mathbf{x}$  must be a nonnull vector. Now if we assume that  $\mathbf{y}$  is positive, then again since  $A$  is positive, we have

$$0 < A\mathbf{y} = A\mathbf{z} - \rho(A)\mathbf{z},$$

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  $i$ . But using Theorem 7.43 implies that  $\rho(A) > \rho(A)$ . Thus,  $\mathbf{y}$  cannot be positive and, since we have already shown that it is nonnegative, we must have  $\mathbf{y} = \mathbf{0}$ . This yields  $A \text{abs}(\mathbf{x}) = \rho(A)\text{abs}(\mathbf{x})$ , so that  $\text{abs}(\mathbf{x})$  is an eigenvector corresponding to  $\rho(A)$ , and from this we get  $\text{abs}(\mathbf{x}) = \rho(A)^{-1}A \text{abs}(\mathbf{x})$ , which shows that  $\text{abs}(\mathbf{x})$  is positive since  $\rho(A) > 0$  and  $A \text{abs}(\mathbf{x}) > \mathbf{0}$ . This completes the proof.  $\square$

An immediate consequence of the proof of Theorem 7.44 is the following.

**Corollary 7.44.1.** Let  $A$  be an  $m \times m$  positive matrix and suppose that  $\lambda$  is an eigenvalue of  $A$  satisfying  $|\lambda| = \rho(A)$ . If  $\mathbf{x}$  is any eigenvector corresponding to  $\lambda$ , then

$$A \text{abs}(\mathbf{x}) = \rho(A)\text{abs}(\mathbf{x})$$

Before determining the dimensionality of the eigenspace associated with the eigenvalue  $\rho(A)$ , we need the following result.

**Theorem 7.45.** Let  $\mathbf{x}$  be an eigenvector corresponding to the eigenvalue  $\lambda$  of the  $m \times m$  positive matrix  $A$ . If  $|\lambda| = \rho(A)$ , then there exists some angle  $\theta$  such that  $e^{-i\theta} \mathbf{x} > \mathbf{0}$ .

*Proof.* Note that

$$\text{abs}(A\mathbf{x}) = \text{abs}(\lambda\mathbf{x}) = \rho(A)\text{abs}(\mathbf{x}), \quad (7.24)$$

while it follows from Corollary 7.44.1 that

$$A \text{abs}(\mathbf{x}) = \rho(A)\text{abs}(\mathbf{x}) \quad (7.25)$$

Now by using (7.24) and (7.25), 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 can happen only if 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, there exists some angle  $\theta$  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  $e^{-i\theta} x_k > 0$  since  $a_{jk} > 0$ .  $\square$

The following result 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 7.46.** 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  $\lambda$  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  $u$  and  $v$  are nonnull vectors satisfying  $Au = \rho(A)u$  and  $Av = \rho(A)v$ , then there exists some scalar  $c$  such that  $v = cu$ . Now from Theorem 7.45, we know there exist angles  $\theta_1$  and  $\theta_2$  such that  $s = e^{-i\theta_1}u > \mathbf{0}$  and  $t = e^{-i\theta_2}v > \mathbf{0}$ . Define  $w = t - ds$ , where

$$d = \min_{1 \leq j \leq m} s_j^{-1} t_j,$$

so that  $w$  is nonnegative with at least one component equal to 0. If  $w \neq \mathbf{0}$ , then clearly  $Aw > \mathbf{0}$  since  $A$  is positive. This leads to a contradiction since

$$Aw = At - dAs = \rho(A)t - \rho(A)ds = \rho(A)w$$

then implies that  $w > \mathbf{0}$ . Thus, we must have  $w = \mathbf{0}$ , so  $t = ds$  and  $v = cu$ , 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  $\lambda$ , of  $A$ . Now if  $x$  is an eigenvector corresponding to  $\lambda$  and  $|\lambda| = \rho(A)$ , then it follows from Theorem 7.45 that there exists an angle  $\theta$  such that  $u = e^{-i\theta}x > \mathbf{0}$ . Clearly,  $Au = \lambda u$ . Premultiplying this identity by  $D_u^{-1}$ , we get

$$D_u^{-1}Au = \lambda \mathbf{1}_m,$$

so that

$$u_i^{-1} \sum_{j=1}^m a_{ij} u_j = \lambda$$

holds for each  $i$ . Now applying Theorem 7.43, we get  $\lambda = \rho(A)$ .  $\square$

We will see that the first statement in the previous theorem actually can be replaced by the stronger condition that  $\rho(A)$  must be a simple eigenvalue of  $A$ . But first we have the following results, the last of which is a very useful limiting result for  $A$ .

**Theorem 7.47.** Suppose that  $A$  is an  $m \times m$  positive matrix, and  $x$  and  $y$  are positive vectors satisfying  $Ax = \rho(A)x$ ,  $A'y = \rho(A)y$ , and  $x'y = 1$ . Then the following hold.

- $(A - \rho(A)xy')^k = A^k - \rho(A)^k xy'$ , for  $k = 1, 2, \dots$
- Each nonzero eigenvalue of  $A - \rho(A)xy'$  is an eigenvalue of  $A$ .
- $\rho(A)$  is not an eigenvalue of  $A - \rho(A)xy'$ .
- $\rho(A - \rho(A)xy') < \rho(A)$ .
- $\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A\}^k = xy'$ .



*Proof.* (a) is easily established by induction, since it clearly holds for  $k = 1$ , and if it holds for  $k = j - 1$ , then

$$\begin{aligned}(A - \rho(A)xy')^j &= (A - \rho(A)xy')^{j-1}(A - \rho(A)xy') \\ &= (A^{j-1} - \rho(A)^{j-1}xy')(A - \rho(A)xy') \\ &= A^j - \rho(A)A^{j-1}xy' - \rho(A)^{j-1}xy'A + \rho(A)^jxy'xy' \\ &= A^j - \rho(A)^jxy' - \rho(A)^jxy' + \rho(A)^jxy' = A^j - \rho(A)^jxy'\end{aligned}$$

Next, suppose that  $\lambda \neq 0$  and  $u$  are an eigenvalue and eigenvector of  $(A - \rho(A)xy')$ , so that

$$(A - \rho(A)xy')u = \lambda u$$

Premultiplying this equation by  $xy'$  and observing that  $xy'(A - \rho(A)xy') = 0$ , we see that we must have  $xy'u = 0$ . Consequently,

$$Au = (A - \rho(A)xy')u = \lambda u,$$

and so  $\lambda$  is also an eigenvalue of  $A$ , as is required for (b). To prove (c), suppose that  $\lambda = \rho(A)$  is an eigenvalue of  $A - \rho(A)xy'$  with  $u$  a corresponding eigenvector. But we have just seen that this would imply that  $u$  is also an eigenvector of  $A$  corresponding to the eigenvalue  $\rho(A)$ . Thus, from Theorem 7.46,  $u = cx$  for some scalar  $c$  and

$$\rho(A)u = (A - \rho(A)xy')u = (A - \rho(A)xy')cx = \rho(A)cx - \rho(A)cx = 0$$

But this is impossible since  $\rho(A) > 0$  and  $u \neq 0$ , and so (c) holds. Now (d) follows directly from (b), (c), and Theorem 7.46. Finally, to prove (e), note that by dividing both sides of the equation given in (a) by  $\rho(A)^k$  and rearranging, we get

$$\{\rho(A)^{-1}A\}^k = xy' + \{\rho(A)^{-1}A - xy'\}^k$$

Take the limit, as  $k \rightarrow \infty$ , of both sides of this equation and observe that from (d),

$$\rho\{\rho(A)^{-1}A - xy'\} = \frac{\rho\{A - \rho(A)xy'\}}{\rho(A)} < 1,$$

and so

$$\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A - xy'\}^k = (0)$$

follows from Theorem 4.23. □

**Theorem 7.48.** 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  $j = r+1, \dots, m$ , due to Theorem 7.46. 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  $r$  1s and  $T_3$  is an upper triangular matrix with each diagonal element equal to 0. Clearly, this limiting matrix has rank at least  $r$ . But from Theorem 7.47(e), we see that the limiting matrix must have rank 1. This proves the result. □

To this point, we have concentrated on positive matrices. Our next step is to extend some of the results above to nonnegative matrices. We will see that many of these results generalize to the class of irreducible nonnegative matrices.

**Definition 7.2.** An  $m \times m$  matrix  $A$ , with  $m \geq 2$ , is called a reducible matrix if there exist some integer  $r$  and  $m \times m$  permutation matrix  $P$  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 the following result regarding irreducible nonnegative matrices.

**Theorem 7.49.** 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  $x$  is an  $m \times 1$  nonnegative vector with  $r$  positive components  $1 \leq r \leq m - 1$ , then  $(I_m + A)x$  has at least  $r + 1$  positive components. Repeated use of this result verifies that  $(I_m + A)^{m-1} > (0)$  since each column of  $I_m + A$  has at least one positive component. Since  $A \geq (0)$ ,  $(I_m + A)x = x + Ax$  must have at least  $r$  positive components. If it has exactly  $r$  positive components, then the  $j$ th component of  $Ax$  must be 0 for every  $j$  for which  $x_j = 0$ . Equivalently, for any permutation matrix  $P$ , the  $j$ th component of  $P Ax$  must be 0 for every  $j$  for which the  $j$ th component of  $P x$  is 0. If we choose a permutation matrix for which  $y = P x$  has its  $m - r$  0s in the last  $m - r$  positions, then we find that the  $j$ th component of  $P Ax = P A P' y$  must be 0 for  $j = r + 1, \dots, m$ . Since  $P A P' \geq (0)$  and the first  $r$  components of  $y$  are positive,  $P A P'$  would have to be of the form

$$P A P' = \begin{bmatrix} B & C \\ (0) & D \end{bmatrix}$$

Since this contradicts the fact that  $A$  is irreducible, the number of positive components in the vector  $(I_m + A)x$  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 since, if for some permutation matrix  $P$ ,

$$P A P' = \begin{bmatrix} B & C \\ (0) & D \end{bmatrix},$$

then

$$P(I_m + A)^{m-1}P' = \begin{bmatrix} I_r + B & C \\ (0) & I_{m-r} + D \end{bmatrix}^{m-1},$$

and the matrix on the right-hand side of this last equation has the upper triangular form given in Definition 7.2.  $\square$

We will generalize the result of Theorem 7.44 by showing that  $\rho(A)$  is pos-

itive, is an eigenvalue of  $A$ , and has a positive eigenvector when  $A$  is an irreducible nonnegative matrix. But first we need the following result.

**Theorem 7.50.** 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 there exists an  $m \times 1$  nonnegative vector  $\mathbf{b}$  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

$$S = \{ \mathbf{y} : \mathbf{y} = (\mathbf{I}_m + A)^{m-1} \mathbf{x}_*, \mathbf{x}_* \in R^m, \mathbf{x}_* \geq \mathbf{0}, \mathbf{x}'_* \mathbf{1}_m = 1 \}$$

Since  $S$  is a closed and bounded set, and  $f$  is a continuous function on  $S$  due to the fact that  $y > 0$  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  $\mathbf{y} = (\mathbf{I}_m + A)^{m-1} \mathbf{x}_*$ . Now it follows from the definition of  $f$  that

$$A \mathbf{x}_* - f(\mathbf{x}_*) \mathbf{x}_* \geq \mathbf{0}$$

Premultiplying this equation by  $(\mathbf{I}_m + A)^{m-1}$  and using the fact that  $(\mathbf{I}_m + A)^{m-1} A = A(\mathbf{I}_m + A)^{m-1}$ , we find that

$$A \mathbf{y} - f(\mathbf{x}_*) \mathbf{y} \geq \mathbf{0}$$

But  $\alpha = f(\mathbf{y})$  is the largest value for which  $A \mathbf{y} - \alpha \mathbf{y} \geq \mathbf{0}$  since 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 7.51.** 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 7.50, 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

$$f(\mathbf{b}) \geq f(m^{-1}\mathbf{1}_m) = \min_{1 \leq i \leq m} (1/m)^{-1}(A)_i(m^{-1}\mathbf{1}_m)$$

$$= \min_{1 \leq i \leq m} \sum_{j=1}^m a_{ij} > 0,$$

since  $A$  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 since  $(I_m + A)^{m-1} > (\mathbf{0})$ , 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  $\mathbf{y} = (I_m + A)^{m-1}\mathbf{b}$ . But  $\alpha = f(\mathbf{y})$  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 since  $\mathbf{b}$  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  $\lambda_i$  is an arbitrary eigenvalue of  $A$ . Now if  $\mathbf{u}$  is an eigenvector of  $A$  corresponding to  $\lambda_i$ , 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

$$A \text{ abs}(\mathbf{u}) - |\lambda_i| \text{ abs}(\mathbf{u}) \geq \mathbf{0},$$

and this 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 7.49, we find that  $\mathbf{b}$  is actually positive. □

The proof of the following result will be left to the reader as an exercise.

**Theorem 7.52.** 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 7.47(e) does not immediately extend to irreducible nonnegative matrices. This leads us to the following definition.

**Definition 7.3.** 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 7.47(e) does extend to primitive matrices and this is summarized below.

**Theorem 7.53.** Let  $A$  be an  $m \times m$  primitive nonnegative matrix and suppose that the  $m \times 1$  vectors  $x$  and  $y$  satisfy  $Ax = \rho(A)x$ ,  $A'y = \rho(A)y$ ,  $x > 0$ ,  $y > 0$ , and  $x'y = 1$ . Then

$$\lim_{k \rightarrow \infty} \{\rho(A)^{-1}A\}^k = xy'$$

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 (1985).

**Theorem 7.54.** Let  $A$  be an  $m \times m$  irreducible nonnegative matrix and suppose that the  $m \times 1$  vectors  $x$  and  $y$  satisfy  $Ax = \rho(A)x$ ,  $A'y = \rho(A)y$ , and  $x'y = 1$ . Then

$$\lim_{N \rightarrow \infty} \left( N^{-1} \sum_{k=1}^N \{\rho(A)^{-1}A\}^k \right) = xy'$$

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 (1990), Medhi (1994), and Taylor and Karlin (1984).

**Example 7.12.** 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,

$$\mathbf{p}^{(1)} = P\mathbf{p}^{(0)}, \quad \mathbf{p}^{(2)} = P\mathbf{p}^{(1)} = P^2\mathbf{p}^{(0)},$$

or for general  $t$ ,

$$\mathbf{p}^{(t)} = P^t\mathbf{p}^{(0)}$$

If we have a large population of individuals subject to the random process discussed above, then  $p_i^{(t)}$  could be described as the proportion of individuals in state  $i$  at time  $t$ , while  $p_i^{(0)}$  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 this depends on the limiting behavior of  $P^t$ , and  $P$  is a nonnegative matrix since each of its elements is a probability. Thus, if  $P$  is a primitive matrix, we can apply Theorem 7.53. Now, since 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 7.43 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 upon 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 his occupation as being upper, middle, or lower class, and these have been labelled 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, or 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 just discussed applies. A simple eigenanalysis of the matrix  $P$  reveals that the positive vector  $\pi$ , satisfying  $P\pi = \pi$  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.

## 10. 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, the reader is referred to Davis (1979).

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

$$A = \begin{bmatrix} a_1 & a_2 & a_3 & \cdots & a_{m-1} & a_m \\ a_m & a_1 & a_2 & \cdots & a_{m-2} & a_{m-1} \\ a_{m-1} & a_m & a_1 & \cdots & a_{m-3} & a_{m-2} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_3 & a_4 & a_5 & \cdots & a_1 & a_2 \\ a_2 & a_3 & a_4 & \cdots & a_m & a_1 \end{bmatrix} \quad (7.26)$$

We will sometimes use the notation  $A = \text{circ}(a_1, a_2, \dots, a_m)$  to refer to the circulant matrix in (7.26). One special circulant matrix, which we will denote by  $\Pi_m$ , 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  $a_1, \dots, a_m$  to denote the columns of an arbitrary  $m \times m$  matrix  $A$  and  $b'_1, \dots, b'_m$  to denote the rows, then

$$A\Pi_m = (a_1, a_2, \dots, a_m)(e_m, e_1, \dots, e_{m-1}) = (a_m, a_1, \dots, a_{m-1}), \tag{7.27}$$

$$\Pi_m A = \begin{bmatrix} e'_2 \\ e'_3 \\ \vdots \\ e'_m \\ e'_1 \end{bmatrix} \begin{bmatrix} b'_1 \\ b'_2 \\ \vdots \\ b'_{m-1} \\ b'_m \end{bmatrix} = \begin{bmatrix} b'_2 \\ b'_3 \\ \vdots \\ b'_m \\ b'_1 \end{bmatrix}, \tag{7.28}$$

and (7.27) equals (7.28) if and only if  $A$  is of the form given in (7.26). Thus we have the following result.

**Theorem 7.55.** The  $m \times m$  matrix  $A$  is a circulant matrix if and only if

$$A = \Pi_m A \Pi'_m$$

Our next theorem gives an expression for an  $m \times m$  circulant matrix in terms of a sum of  $m$  matrices.

**Theorem 7.56.** 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 (7.26), we see that

$$A = a_1 I_m + a_2(e_m, e_1, \dots, e_{m-1}) + a_3(e_{m-1}, e_m, e_1, \dots, e_{m-2}) + \dots + a_m(e_2, e_3, \dots, e_m, e_1)$$

Since the postmultiplication of any  $m \times m$  matrix by  $\Pi_m$  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 the following theorem.

**Theorem 7.57.** Let  $A$  and  $B$  be  $m \times m$  circulant matrices. Then

- (a)  $A'$  is circulant,
- (b) for any scalars  $\alpha$  and  $\beta$ ,  $\alpha A + \beta B$  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 (7.26) 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  $A = \Pi_m A \Pi_m'$ . But  $\Pi_m$  is an orthogonal matrix, so

$$A^r = (\Pi_m A \Pi_m')^r = \Pi_m A^r \Pi_m',$$

and so by Theorem 7.55,  $A^r$  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 7.56 provides a simple way of proving the following result.

**Theorem 7.58.** 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 7.56 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,

$$\begin{aligned} AB &= \left( \sum_{i=1}^m a_i \Pi_m^{i-1} \right) \left( \sum_{j=1}^m b_j \Pi_m^{j-1} \right) = \sum_{i=1}^m \sum_{j=1}^m (a_i \Pi_m^{i-1})(b_j \Pi_m^{j-1}) \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i b_j \Pi_m^{i+j-2} = \sum_{i=1}^m \sum_{j=1}^m (b_j \Pi_m^{j-1})(a_i \Pi_m^{i-1}) \\ &= \left( \sum_{j=1}^m b_j \Pi_m^{j-1} \right) \left( \sum_{i=1}^m a_i \Pi_m^{i-1} \right) = BA \quad \square \end{aligned}$$

All circulant matrices are diagonalizable. We will show this by determining the eigenvalues and eigenvectors of a circulant matrix. But first let us find the eigenvalues and eigenvectors of the special circulant matrix  $\Pi_m$ .

**Theorem 7.59.** 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  $\Lambda$  to be the diagonal matrix  $\text{diag}(1, \theta, \dots, \theta^{m-1})$  and let

$$F = \frac{1}{\sqrt{m}} \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & \theta & \theta^2 & \dots & \theta^{m-1} \\ 1 & \theta^2 & \theta^4 & \dots & \theta^{2(m-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \theta^{m-1} & \theta^{2(m-1)} & \dots & \theta^{(m-1)(m-1)} \end{bmatrix}$$

Then the diagonalization of  $\Pi_m$  is given by  $\Pi_m = F\Lambda F^*$ , where  $F^*$  is the

conjugate transpose of  $F$ ; that is, the diagonal elements of  $\Lambda$  are the eigenvalues of  $\Pi_m$ , while the columns of  $F$  are corresponding eigenvectors.

*Proof.* The eigenvalue–eigenvector equation,  $\Pi_m x = \lambda x$ , 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  $\Pi_m$  are  $1, \theta, \dots, \theta^{m-1}$ . Substituting the eigenvalue  $\theta^{j-1}$  and  $x_1 = m^{-1/2}$  into the equations above, we find that an eigenvector corresponding to the eigenvalue  $\theta^{j-1}$  is given by  $x = m^{-1/2}(1, \theta^{j-1}, \dots, \theta^{(m-1)(j-1)})'$ . Thus, we have shown that the diagonal elements of  $\Lambda$  are the eigenvalues of  $\Pi_m$  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 7.59 is sometimes referred to as the Fourier matrix of order  $m$ . The diagonalization of an arbitrary circulant matrix, which follows directly from Theorems 7.56 and 7.59, is given in our next theorem.

**Theorem 7.60.** 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  $\lambda_j$  and  $F$  are defined as in Theorem 7.59.

*Proof.* Since  $\Pi_m = F\Lambda F^*$  and  $FF^* = I_m$ , we have  $\Pi_m^j = F\Lambda^j F^*$ , for  $j = 2, \dots, m-1$ , and so by using Theorem 7.56, we find that

$$\begin{aligned} A &= a_1 I_m + a_2 \Pi_m + a_3 \Pi_m^2 + \dots + a_m \Pi_m^{m-1} \\ &= a_1 FF^* + a_2 F\Lambda^1 F^* + a_3 F\Lambda^2 F^* + \dots + a_m F\Lambda^{m-1} F^* \\ &= F(a_1 I_m + a_2 \Lambda^1 + a_3 \Lambda^2 + \dots + a_m \Lambda^{m-1})F^* = F\Delta F^* \quad \square \end{aligned}$$

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

$$A = \begin{bmatrix} a_0 & a_1 & a_2 & \cdots & a_{m-2} & a_{m-1} \\ a_{-1} & a_0 & a_1 & \cdots & a_{m-3} & a_{m-2} \\ a_{-2} & a_{-1} & a_0 & \cdots & a_{m-4} & a_{m-3} \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ a_{-m+2} & a_{-m+3} & a_{-m+4} & \cdots & a_0 & a_1 \\ a_{-m+1} & a_{-m+2} & a_{-m+3} & \cdots & a_{-1} & a_0 \end{bmatrix}$$

If  $a_j = a_{-j}$  for  $j = 1, \dots, m - 1$ , then the matrix  $A$  is a symmetric Toeplitz matrix. One important and fairly simple symmetric Toeplitz matrix is one that has  $a_j = a_{-j} = 0$  for  $j = 2, \dots, m - 1$ , so that

$$A = \begin{bmatrix} a_0 & a_1 & 0 & \cdots & 0 & 0 \\ a_1 & a_0 & a_1 & \cdots & 0 & 0 \\ 0 & a_1 & a_0 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & a_0 & a_1 \\ 0 & 0 & 0 & \cdots & a_1 & a_0 \end{bmatrix} \tag{7.29}$$

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).

### 11. 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

$$H'H = HH' = mI_m; \tag{7.30}$$

that is, the columns of  $H$  form an orthogonal set of vectors, and the rows form an orthogonal set as well. For instance, a  $2 \times 2$  Hadamard matrix is given by

$$H = \begin{bmatrix} 1 & 1 \\ 1 & -1 \end{bmatrix},$$

while a  $4 \times 4$  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 the following theorem.

**Theorem 7.61.** 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 (7.30). Also using (7.30), we find that

$$|H'_m H_m| = |mI_m| = m^m$$

But

$$|H'_m H_m| = |H'_m| |H_m| = |H_m|^2,$$

and so (b) follows. To prove (c), note that each element of  $H_m \otimes H_n$  is +1 or -1 since each element is the product of an element from  $H_m$  and an element from  $H_n$ , and

$$(H_m \otimes H_n)'(H_m \otimes H_n) = H'_m H_m \otimes H'_n H_n = mI_m \otimes nI_n = mnI_{mn} \quad \square$$

Hadamard matrices which 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 7.62.** If there exists an  $m \times m$  Hadamard matrix, then there exists an  $m \times m$  normalized Hadamard matrix.

*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  $D^2 = I_m$  since 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,

$$H'_* H_* = (HD)' HD = D' H' HD = D(mI_m)D = mD^2 = mI_m$$

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  $2 \times 2$  Hadamard matrix, and this matrix can be used repeatedly in Theorem 7.61(c) to obtain a  $2^n \times 2^n$  Hadamard matrix for any  $n \geq 2$ . However,  $m \times m$  Hadamard matrices do exist for some values of  $m \neq 2^n$ . Our next result gives a necessary condition on the order  $m$  so that Hadamard matrices of order  $m$  exist.

**Theorem 7.63.** If  $H$  is an  $m \times m$  Hadamard matrix, where  $m > 2$ , then  $m$  is a multiple of 4.

*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, due to Theorem 7.62, 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  $r$  +1s and  $r$  -1s, where  $r = n/2$ ; thus clearly,

$$n = 2r, \tag{7.31}$$

or in other words,  $n$  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 since  $n_{++} + n_{+-} = r$ ,  $n_{++} + n_{-+} = r$ , and  $n_{--} + n_{+-} = r$ . For instance, if  $n_{++} = s$ , then  $n_{+-} = (r - s)$ ,  $n_{-+} = (r - s)$ , and  $n_{--} = s$ . But 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 (7.31) we get  $n = 4s$ , which completes the proof.  $\square$

Some additional results on Hadamard matrices can be found in Hedayat and Wallis (1978) and Agaian (1985).

An  $m \times m$  matrix  $A$  is said to be a Vandermonde matrix if it has the form

$$A = \begin{bmatrix} 1 & 1 & 1 & \cdots & 1 \\ a_1 & a_2 & a_3 & \cdots & a_m \\ a_1^2 & a_2^2 & a_3^2 & \cdots & a_m^2 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a_1^{m-1} & a_2^{m-1} & a_3^{m-1} & \cdots & a_m^{m-1} \end{bmatrix} \quad (7.32)$$

For instance, if  $F$  is the  $m \times m$  Fourier matrix discussed in Section 7.10 then  $A = m^{1/2}F$  is a Vandermonde matrix with  $a_i = \theta^{i-1}$ , for  $i = 1, \dots, m$ . Our final result of this chapter gives an expression for the determinant of a Vandermonde matrix.

**Theorem 7.64.** Let  $A$  be the  $m \times m$  Vandermonde matrix given in (7.31). Then its determinant is given by

$$|A| = \prod_{1 \leq i < j \leq m} (a_j - a_i) \quad (7.33)$$

*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 (7.33) holds when  $A$  is  $2 \times 2$ . Next we assume that (7.33) holds for Vandermonde matrices of order  $m - 1$  and show that then 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, since  $B$  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

$$C = \begin{bmatrix} 1 & 0 & 0 & \cdots & 0 & 0 \\ -a_1 & 1 & 0 & \cdots & 0 & 0 \\ 0 & -a_1 & 1 & \cdots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \cdots & 1 & 0 \\ 0 & 0 & 0 & \cdots & -a_1 & 1 \end{bmatrix},$$



## PROBLEMS

and note that by repeatedly using the cofactor expansion formula for a determinant on the first row, we find that  $|C| = 1$ . Thus,  $|A| = |CA|$ . But 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,

$$\begin{aligned} |A| = |E| = |BD| = |B||D| &= \left\{ \prod_{2 \leq i < j \leq m} (a_j - a_i) \right\} \left\{ \prod_{2 \leq j \leq m} (a_j - a_1) \right\} \\ &= \prod_{1 \leq i < j \leq m} (a_j - a_i), \end{aligned}$$

where the second equality was obtained by using the cofactor expansion formula on the first column of  $E$ . This completes the proof.  $\square$

## PROBLEMS

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.

- Give an expression for the determinant of  $A$ .
- For what values of  $a$ ,  $b$ ,  $c$ , and  $d$  will  $A$  be nonsingular?
- Find an expression for  $A^{-1}$ .

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 utilizing equations (7.2)–(7.5).

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.

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  $G$ .

5. Use Theorems 7.1 and 7.4 to find the determinant and inverse of the matrix

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

6. Let  $A$  be an  $m \times n$  matrix partitioned as

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

where  $A_{11}$  is  $r \times r$  and  $\text{rank}(A) = \text{rank}(A_{11}) = r$ .

- (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

$$A^+ = \begin{bmatrix} A'_{11} \\ A'_{12} \end{bmatrix} C [A'_{11} \quad A'_{21}],$$

where  $C = (A_{11}A'_{11} + A_{12}A'_{12})^{-1} A_{11} (A'_{11}A_{11} + A'_{21}A_{21})^{-1}$ .

7. Use Theorem 7.4 to show that if

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

is nonsingular and  $A_{11}$  is nonsingular, then  $A_{22} - A_{21}A_{11}^{-1}A_{12}$  is nonsingular.

8. 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  $r \times r$  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  $m - r$ .

9. Consider the  $m \times m$  matrix

$$A = \begin{bmatrix} A_{11} & a \\ a' & a_{mm} \end{bmatrix},$$

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  $a = 0$ .

(b) Use part (a) to obtain an alternative proof of Theorem 7.23; that is, 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.

10. Let  $A$  be an  $m \times m$  matrix and define  $A_i$  to be the  $i \times i$  matrix obtained by deleting the last  $m - i$  rows and columns of  $A$ . The leading principal minors of  $A$  are given by the determinants,  $|A_1|, \dots, |A_m|$ , where  $A_m = A$ .

Show that if  $A$  is a symmetric matrix, then it is positive definite if and only if all of its leading principal minors are positive.

11. Let the  $2 \times 2$  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  $(A \otimes B)^{-1}$ .

12. Give a simplified expression for  $I_m \otimes I_n$ .

13. Prove the properties given in Theorem 7.6.

14. Prove results (b) and (c) of Theorem 7.9.

15. Show that if  $A$  and  $B$  are symmetric matrices, then  $A \otimes B$  is also symmetric.

16. Find the rank of  $A \otimes B$ , where

$$A = \begin{bmatrix} 2 & 6 \\ 1 & 4 \\ 3 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 5 & 2 & 4 \\ 2 & 1 & 1 \\ 1 & 0 & 2 \end{bmatrix}$$

17. For matrices  $A$  and  $B$  of any size, show that  $A \otimes B = (0)$  if and only if  $A = (0)$  or  $B = (0)$ .

18. Let  $x_i$  be an eigenvector of the  $m \times m$  matrix  $A$  corresponding to the eigenvalue  $\lambda_i$ . Let  $y_j$  be an eigenvector of the  $p \times p$  matrix  $B$  corresponding to the eigenvalue  $\theta_j$ .

- (a) Show that  $x_i \otimes y_j$  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  $B$ .

19. Show that if  $A$  and  $B$  are positive definite matrices, then  $A \otimes B$  is also positive definite.

## PROBLEMS

20. Let  $x$  be an  $m \times 1$  vector and  $y$  be an  $n \times 1$  vector. Verify that the three matrices  $xy'$ ,  $y' \otimes x$ , and  $x \otimes y'$  are identical.
21. Compute the sum of squared errors  $SSE = (y - \hat{y})'(y - \hat{y})$  for the two-way classification model with interaction discussed in Example 7.5.
22. 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  $\beta = (\mu, \tau_1, \dots, \tau_a, \gamma_1, \dots, \gamma_b)'$ , 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

$$SSA = nb \sum_{i=1}^a (\bar{y}_{i.} - \bar{y}_{..})^2$$

- (c) In a similar fashion, show that the sum of squares for factor B is

$$SSB = na \sum_{j=1}^b (\bar{y}_{.j} - \bar{y}_{..})^2$$

- (d) Find a set of as many linearly independent estimable functions of  $\mu$ ,  $\tau_i$ , and  $\gamma_j$  as possible.
- (e) Use the sum of squared errors computed in (a) and the sum of squared errors computed in Problem 21 to show that the sum of squares for interaction in the model of Problem 21 is given by

$$SSAB = n \sum_{i=1}^a \sum_{j=1}^b (\bar{y}_{ij} - \bar{y}_{i.} - \bar{y}_{.j} + \bar{y}_{..})^2$$

23. Prove Theorem 7.13.

24. 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)  $(A_1 \oplus A_2)(A_3 \oplus A_4) = A_1A_3 \oplus A_2A_4$ ,  
 (c)  $(A_1 \oplus A_2) \otimes A_3 = (A_1 \otimes A_3) \oplus (A_2 \otimes A_3)$ .

25. Give an example to show that, in general,

$$A_1 \otimes (A_2 \oplus A_3) \neq (A_1 \otimes A_2) \oplus (A_1 \otimes A_3)$$

26. Complete the details of Example 7.7; that is, use Theorem 6.4 and Theorem 7.16 to prove Theorem 6.5.

27. Prove the results of Corollary 7.17.1.

28. Let  $A$  and  $B$  be  $m \times n$  and  $n \times p$  matrices, respectively, while  $\mathbf{c}$  and  $\mathbf{d}$  are  $p \times 1$  and  $n \times 1$  vectors. Show that

(a)  $AB\mathbf{c} = (\mathbf{c}' \otimes A) \text{vec}(B) = (A \otimes \mathbf{c}') \text{vec}(B')$ ,

(b)  $\mathbf{d}'B\mathbf{c} = (\mathbf{c}' \otimes \mathbf{d}') \text{vec}(B)$ .

29. For any matrix  $A$  and any vector  $\mathbf{b}$ , show that

$$\text{vec}(A \otimes \mathbf{b}) = \text{vec}(A) \otimes \mathbf{b}$$

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

$$\text{vec}(AC + CB) = \{(I_n \otimes A) + (B' \otimes I_m)\} \text{vec}(C)$$

31. If  $\mathbf{e}_i$  is the  $i$ th column of the identity matrix  $I_m$ , verify that

$$\text{vec}(I_m) = \sum_{i=1}^m (\mathbf{e}_i \otimes \mathbf{e}_i)$$

32. Prove property (h) of Theorem 7.18.

33. Let the  $2 \times 2$  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 7.22?

34. Give an example of matrices  $A$  and  $B$  such that neither is nonnegative definite, yet  $A \odot B$  is positive definite.
35. Let  $A$ ,  $B$ , and  $C$  be  $m \times n$  matrices. Show that

$$\operatorname{tr}\{(A' \odot B')C\} = \operatorname{tr}\{A'(B \odot C)\}$$

36. Suppose that the  $m \times m$  matrix  $A$  is diagonalizable; that is, there exist a nonsingular matrix  $X$  and a diagonal matrix  $\Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_m)$  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

$$(X \odot X'^{-1})\mathbf{1}_m = (X \odot X'^{-1})'\mathbf{1}_m = \mathbf{1}_m$$

37. Let  $A$  and  $B$  be  $m \times m$  nonnegative definite matrices. Show that
- $|A \odot B| \geq |A||B|$ ,
  - $|A \odot A^{-1}| \geq 1$ , if  $A$  is positive definite.
38. For each of the following pairs of  $2 \times 2$  matrices, compute the smaller eigenvalue  $\lambda_2(A \odot B)$  and the lower bounds for this eigenvalue given by Theorem 7.26 and Theorem 7.28. Which bound is closer to the actual value?

$$(a) \quad A = \begin{bmatrix} 4 & 0 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & 0 \\ 0 & 3 \end{bmatrix}$$

$$(b) \quad A = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}, \quad B = \begin{bmatrix} 2 & \sqrt{2} \\ \sqrt{2} & 3 \end{bmatrix}$$

39. Let  $A$  be an  $m \times m$  positive definite matrix. Use Theorem 7.24 to show that, if  $B = A^{-1}$  then  $a_{ii}b_{ii} \geq 1$ . Show how this generalizes to  $a_{ii}b_{ii} \geq 1$  for  $i = 1, \dots, m$ .
40. 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 7.25 on  $|R_A \odot C|$ , where  $C = R_B - (e_1' R_B^{-1} e_1)^{-1} e_1 e_1'$ , to establish the inequality given in (a).

41. 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.
42. 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  $\text{rank}(A \odot B) = r$ .
43. Show that if  $A$  and  $B$  are singular  $2 \times 2$  matrices then  $A \odot B$  is also singular.
44. Let  $R$  be an  $m \times m$  positive definite correlation matrix having  $\lambda$  as its smallest eigenvalue. Show that if  $\tau$  is the smallest eigenvalue of  $R \odot R$  and  $R \neq I_m$ , then  $\tau > \lambda$ .
45. Consider the matrix

$$\Psi_m = \sum_{i=1}^m e_{i,m} (e_{i,m} \otimes e_{i,m})',$$

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  $w(A)$  to be the  $m \times 1$  vector containing the diagonal elements of  $A$ ; that is,  $w(A) = (a_{11}, \dots, a_{mm})'$ . Also let  $\Lambda_m$  be the  $m^2 \times m^2$  matrix given by

$$\Lambda_m = \sum_{i=1}^m (E_{ii} \otimes E_{ii}) = \sum_{i=1}^m (e_{i,m} e_{i,m}' \otimes e_{i,m} e_{i,m}')$$

Show that

- (a)  $\Psi_m' w(A) = \text{vec}(A)$  for every diagonal matrix  $A$ ,
- (b)  $\Psi_m \text{vec}(A) = w(A)$  for every matrix  $A$ ,
- (c)  $\Psi_m \Psi_m' = I_m$  so that  $\Psi_m^+ = \Psi_m'$ ,
- (d)  $\Psi_m' \Psi_m = \Lambda_m$ ,



## PROBLEMS

$$(e) \Lambda_m N_m = N_m \Lambda_m = \Lambda_m,$$

$$(f) \{\text{vec}(A)\}' \Lambda_m (B \otimes B) \Lambda_m \text{vec}(A) = \{w(A)\}' (B \odot B) w(A).$$

Additional properties of  $\Psi_m$  can be found in Magnus (1988).

46. 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  $K_{mn}$ .

47. Write out the commutation matrices  $K_{22}$  and  $K_{24}$ .

48. The eigenvalues of  $K_{mm}$  were given in Theorem 7.33. Show that corresponding eigenvectors are given by the vectors of the form  $e_l \otimes e_l$ ,  $(e_l \otimes e_k) + (e_k \otimes e_l)$ , and  $(e_l \otimes e_k) - (e_k \otimes e_l)$ .

49. Show that the commutation matrix  $K_{mn}$  can be expressed as

$$K_{mn} = \sum_{i=1}^m (e_i \otimes I_n \otimes e_i'),$$

where  $e_i$  is the  $i$ th column of  $I_m$ . Use this to show that if  $A$  is  $n \times m$ ,  $x$  is  $m \times 1$ ,  $y$  is an arbitrary vector, then

$$K'_{mn} (x \otimes A \otimes y') = A \otimes xy'$$

50. 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)  $P$  is symmetric,

(b)  $\text{rank}(P) = r^2$ ,

(c)  $\text{tr}(P) = \text{tr}(A'A)$ ,

(d)  $P^2 = (AA') \otimes (A'A)$ ,

(e) the nonzero eigenvalues of  $P$  are  $\lambda_1, \dots, \lambda_r$  and  $\pm(\lambda_i \lambda_j)^{1/2}$  for all  $i < j$ .

51. Prove the results of Theorem 7.35.

52. Show that if  $A$  and  $B$  are  $m \times m$  matrices, then

$$\begin{aligned} N_m(A \otimes B + B \otimes A)N_m &= (A \otimes B + B \otimes A)N_m = N_m(A \otimes B + B \otimes A) \\ &= 2N_m(A \otimes B)N_m \end{aligned}$$

53. Write out the matrices  $N_2$  and  $N_3$ .

54. For  $i = 1, \dots, m, j = 1, \dots, i$ , define the  $m(m+1)/2 \times 1$  vector  $u_{ij}$  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,

$$I_{m(m+1)/2} = (u_{11}, u_{21}, \dots, u_{m1}, u_{22}, \dots, u_{m2}, u_{33}, \dots, u_{mm})$$

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  $D_m = \sum_{i \geq j} \{\text{vec}(T_{ij})\} u'_{ij}$ ; that is, verify that

$$\sum_{i \geq j} \{\text{vec}(T_{ij})\} u'_{ij} v(A) = \text{vec}(A),$$

where  $A$  is an arbitrary  $m \times m$  symmetric matrix.

55. Prove the results of Theorem 7.40.

56. If  $A$  is an  $m \times m$  matrix show that

$$(a) \quad D_m D_m^+ (A \otimes A) D_m = (A \otimes A) D_m,$$

$$(b) \quad \{D_m^+ (A \otimes A) D_m\}^i = D_m^+ (A^i \otimes A^i) D_m, \text{ where } i \text{ is any positive integer.}$$

57. If  $u_{ij}$  and  $E_{ij}$  are defined as in Problem 54, show that  $L'_m = \sum_{i \geq j} \{\text{vec}(E_{ij})\} u'_{ij}$ ; that is, verify that

$$\sum_{i \geq j} \{\text{vec}(E_{ij})\} u'_{ij} v(A) = \text{vec}(A),$$

where  $A$  is an arbitrary  $m \times m$  lower triangular matrix.

58. Prove Theorem 7.41.

59. For  $i = 2, \dots, m, j = 1, \dots, i - 1$ , define the  $m(m - 1)/2 \times 1$  vector  $\tilde{u}_{ij}$  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,

$$I_{m(m-1)/2} = (\tilde{u}_{21}, \dots, \tilde{u}_{m1}, \tilde{u}_{32}, \dots, \tilde{u}_{m2}, \tilde{u}_{43}, \dots, \tilde{u}_{mm-1})$$

Show that  $\tilde{L}'_m = \sum_{i>j} \{\text{vec}(E_{ij})\} \tilde{u}'_{ij}$ ; that is, verify that

$$\sum_{i>j} \{\text{vec}(E_{ij})\} \tilde{u}'_{ij} \tilde{v}(A) = \text{vec}(A),$$

where  $A$  is an arbitrary  $m \times m$  strictly lower triangular matrix.

60. Prove the results of Theorem 7.42.
61. Find a  $2 \times 2$  nonnegative matrix  $A$  which has its spectral radius equal to 1, yet  $A^k$  does not converge to anything as  $k \rightarrow \infty$ .
62. Show that if  $A$  is a nonnegative matrix and, for some positive integer  $k, A^k$  is a positive matrix, then  $\rho(A) > 0$ .
63. It can be shown [see, for example, Horn and Johnson (1985)] that if  $A$  is an  $m \times m$  nonnegative matrix, then  $\rho(A)$  is an eigenvalue of  $A$  and there exists a nonnegative eigenvector  $x$  corresponding to the eigenvalue  $\rho(A)$ . This result is weaker than the result for irreducible nonnegative matrices. For each of the following, find a  $2 \times 2$  nonnull reducible matrix  $A$  such that the stated condition holds.
- $\rho(A) = 0$ .
  - $x$  is not positive for any  $x$  satisfying  $Ax = \rho(A)x$ .
  - $\rho(A)$  is a multiple eigenvalue.
64. Verify that the absolute value of each of the eigenvalues of the  $2 \times 2$  irreducible matrix

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

is equal to  $\rho(A)$ .

65. Let  $A$  be an  $m \times m$  irreducible nonnegative matrix.
- 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 7.52; that is, prove that for any irreducible nonnegative matrix  $A$ ,  $\rho(A)$  must be a simple eigenvalue.

66. 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  $\pi$  such that  $\lim_{t \rightarrow \infty} p^{(t)} = \pi$ .
67. 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  $A$ .

68. Show that the conjugate transpose of the matrix  $F$  given in Theorem 7.59 is

$$F^* = \frac{1}{\sqrt{m}} \begin{bmatrix} 1 & 1 & 1 & \dots & 1 \\ 1 & \theta^{-1} & \theta^{-2} & \dots & \theta^{-(m-1)} \\ 1 & \theta^{-2} & \theta^{-4} & \dots & \theta^{-2(m-1)} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \theta^{-(m-1)} & \theta^{-2(m-1)} & \dots & \theta^{-(m-1)(m-1)} \end{bmatrix}$$

Then use the geometric series partial sum formula

$$\sum_{j=0}^n r^j = \frac{1 - r^{n+1}}{1 - r}$$

to prove that  $F^{-1} = F^*$ .

69. Let  $F$  be defined as in Theorem 7.59 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)  $F^3 = F^*$ .

## PROBLEMS

70. Let  $\Pi_m$  be the circulant matrix defined in Section 7.10. Show that
- $\Pi_m^{m-1} = \Pi_m^{-1}$ ,
  - $\Pi_m^m = I_m$ ,
  - $\Pi_m^{mn+r} = \Pi_m^r$ , for any integers  $n$  and  $r$ .
71. 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  $AB$ .
72. Use Theorem 7.60 to find the eigenvalues of the circulant matrix  $A = \text{circ}(1, \dots, 1)$ .
73. Show that if  $A$  is a singular circulant matrix, then its Moore–Penrose inverse,  $A^+$ , is also a circulant matrix.
74. 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.
75. 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

$$A = a_0 I_m + \sum_{j=1}^{m-1} (a_j B^j + a_{-j} B'^j)$$

76. Consider the  $m \times m$  Toeplitz matrix

$$A = \begin{bmatrix} 1 & b & b^2 & \dots & b^{m-1} \\ a & 1 & b & \dots & b^{m-2} \\ a^2 & a & 1 & \dots & b^{m-3} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ a^{m-1} & a^{m-2} & a^{m-3} & \dots & 1 \end{bmatrix},$$

where  $ab \neq 1$ . Verify by multiplication that the inverse of  $A$  is given by

$$A^{-1} = \begin{bmatrix} c & -bc & 0 & \dots & 0 & 0 \\ -ac & (ab+1)c & -bc & \dots & 0 & 0 \\ 0 & -ac & (ab+1)c & \dots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & 0 & \dots & (ab+1)c & -bc \\ 0 & 0 & 0 & \dots & -ac & c \end{bmatrix},$$

where  $c = (1 - ab)^{-1}$ . Show that  $A$  is singular if  $ab = 1$ .

77. Suppose that  $z_1, \dots, z_{m+1}$  are independent random variables each having mean 0 and variance 1. Let  $x$  be the  $m \times 1$  random vector that has as its  $i$ th component

$$x_i = z_{i+1} - \rho z_i,$$

where  $\rho$  is a constant. Show that the covariance matrix of  $x$  is a Toeplitz matrix of the form given in (7.29), and find the values of  $a_0$  and  $a_1$ .

78. Find a Hadamard matrix of order 8.
79. 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  $n$ .
80. Show that the determinant of a Hadamard matrix attains the upper bound of the Hadamard inequality given in Corollary 7.23.1.
81. 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  $4m$ .

82. Show that the Vandermonde matrix  $A$  given in (7.32) is nonsingular if and only if the  $m$  elements of the second row are distinct.
83. Let  $A$  be the  $m \times m$  Vandermonde matrix given in (7.32). Prove that if there are  $r$  distinct values in the set  $\{a_1, \dots, a_m\}$ , then  $\text{rank}(A) = r$ .
84. 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.

## CHAPTER EIGHT

# Matrix Derivatives and Related Topics

### 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 utilize 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.

### 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 (1988). 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,

$$f(x+u) = f(x) + uf'(x) + r_1(u, x), \quad (8.1)$$

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

$$d_u f(x) = u f'(x) \quad (8.2)$$

appearing in (8.1) 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 (8.2) simply by  $df$ . Generalizations of (8.1) can be obtained by taking higher-order 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

$$y'(x) = g'(f(x)) f'(x) \quad (8.3)$$

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

$$\frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) = \left[ \frac{\partial}{\partial x_1} f(\mathbf{x}) \cdots \frac{\partial}{\partial x_n} f(\mathbf{x}) \right],$$

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  $\mathbf{I}_n$ . The first-order Taylor formula analogous to (8.1) is given by

$$f(\mathbf{x} + \mathbf{u}) = f(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) \right) \mathbf{u} + r_1(\mathbf{u}, \mathbf{x}), \quad (8.4)$$

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 (8.4) 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 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,  $d^2f$ , can be written as a quadratic form in the vector  $\mathbf{u}$ ; that is,

$$d^2f = \mathbf{u}' H_f \mathbf{u},$$

where  $H_f$ , called the Hessian matrix, is the matrix of second-order partial derivatives given by

$$H_f = \begin{bmatrix} \frac{\partial^2}{\partial x_1^2} f(\mathbf{x}) & \frac{\partial^2}{\partial x_1 \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_1 \partial x_n} f(\mathbf{x}) \\ \frac{\partial^2}{\partial x_2 \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_2^2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_2 \partial x_n} f(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2}{\partial x_n \partial x_1} f(\mathbf{x}) & \frac{\partial^2}{\partial x_n \partial x_2} f(\mathbf{x}) & \cdots & \frac{\partial^2}{\partial x_n^2} f(\mathbf{x}) \end{bmatrix}$$

### 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 componentwise to  $\mathbf{f}$ . For instance, the first-order Taylor formula is given by

$$\mathbf{f}(\mathbf{x} + \mathbf{u}) = \mathbf{f}(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}'} \mathbf{f}(\mathbf{x}) \right) \mathbf{u} + r_1(\mathbf{u}, \mathbf{x}) = \mathbf{f}(\mathbf{x}) + d\mathbf{f}(\mathbf{x}) + r_1(\mathbf{u}, \mathbf{x}),$$

where the vector remainder,  $r_1(u, x)$ , satisfies

$$\lim_{u \rightarrow 0} \frac{r_1(u, x)}{(u'u)^{1/2}} = 0$$

and the first derivative of  $f$  at  $x$  is given by

$$\frac{\partial}{\partial x'} f(x) = \begin{bmatrix} \frac{\partial}{\partial x_1} f_1(x) & \frac{\partial}{\partial x_2} f_1(x) & \cdots & \frac{\partial}{\partial x_n} f_1(x) \\ \frac{\partial}{\partial x_1} f_2(x) & \frac{\partial}{\partial x_2} f_2(x) & \cdots & \frac{\partial}{\partial x_n} f_2(x) \\ \vdots & \vdots & & \vdots \\ \frac{\partial}{\partial x_1} f_m(x) & \frac{\partial}{\partial x_2} f_m(x) & \cdots & \frac{\partial}{\partial x_n} f_m(x) \end{bmatrix}$$

This matrix of partial derivatives is sometimes referred to as the Jacobian matrix of  $f$  at  $x$ . Again, it is crucial to understand the relationship between the first differential and the first derivative. If we obtain the first differential of  $f$  at  $x$  in  $u$  and write it in the form

$$df = Bu,$$

then the  $m \times n$  matrix  $B$  must be the derivative of  $f$  at  $x$ .

If  $y$  and  $g$  are real-valued functions satisfying  $y(x) = g(f(x))$ , then the generalization of the chain rule given in (8.3) is

$$\frac{\partial}{\partial x_i} y(x) = \sum_{j=1}^m \left( \frac{\partial}{\partial f_j} g(f) \right) \left( \frac{\partial}{\partial x_i} f_j(x) \right) = \left( \frac{\partial}{\partial f'} g(f) \right) \left( \frac{\partial}{\partial x_i} f(x) \right)$$

for  $i = 1, \dots, n$ , or simply

$$\frac{\partial}{\partial x'} y(x) = \left( \frac{\partial}{\partial f'} g(f) \right) \left( \frac{\partial}{\partial x'} f(x) \right)$$

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  $f(x)$  can be easily extended to the matrix function  $F(X)$  by utilizing the  $\text{vec}$  operator; that is, let  $f$  be the  $pq \times 1$  vector function such that  $f(\text{vec}(X)) = \text{vec}(F(X))$ . Then, for instance, the Jacobian matrix of  $F$  at  $X$  is given by the  $pq \times mn$  matrix

$$\frac{\partial}{\partial \text{vec}(X)'} f(\text{vec}(X)) = \frac{\partial}{\partial \text{vec}(X)'} \text{vec}(F(X)),$$

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

$$\text{vec}(d^i F) = \text{vec}(d_U^i F(X)) = d^i f = d_{\text{vec}(U)}^i f(\text{vec}(X));$$

that is,  $d^i F$ , the  $i$ th order differential of  $F$  at  $X$  in the incremental matrix  $U$ , is defined to be the  $p \times q$  matrix obtained by unstacking the  $i$ th-order differential of  $f$  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  $\alpha$  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 the properties above 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)  $d(XY) = (dX)Y + X(dY)$ .

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 use of some of these properties first by finding the derivatives of some simple scalar functions of a vector  $x$ , and then by finding the derivatives of some simple matrix functions of a matrix  $X$ .

**Example 8.1.** Let  $x$  be an  $m \times 1$  vector of unrelated variables and define the functions

$$f(x) = a'x,$$

where  $a$  is an  $m \times 1$  vector of constants, and

$$g(x) = x'Ax,$$

where  $A$  is an  $m \times m$  symmetric matrix of constants. 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 and derivative of our second function are given by

$$\begin{aligned} dg &= d(\mathbf{x}'A\mathbf{x}) = d(\mathbf{x}')A\mathbf{x} + \mathbf{x}' d(A\mathbf{x}) = (d\mathbf{x})'A\mathbf{x} + \mathbf{x}'A d\mathbf{x} \\ &= \{(d\mathbf{x})'A\mathbf{x}\}' + \mathbf{x}'A d\mathbf{x} = \mathbf{x}'A' d\mathbf{x} + \mathbf{x}'A d\mathbf{x} = 2\mathbf{x}'A d\mathbf{x}, \end{aligned}$$

and

$$\frac{\partial}{\partial \mathbf{x}'} g = 2\mathbf{x}'A$$

**Example 8.2.** 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. For our first function, we find that

$$dF = d(AX) = A dX,$$

so that

$$d \operatorname{vec}(F) = \operatorname{vec}(dF) = \operatorname{vec}(A dX) = (\mathbf{I}_n \otimes A) \operatorname{vec}(dX) = (\mathbf{I}_n \otimes A) d \operatorname{vec}(X)$$

Thus, we must have

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(F) = I_n \otimes A$$

The differential of our second function is

$$\begin{aligned} dG &= d\{(X - C)'B(X - C)\} \\ &= \{d(X' - C')\}B(X - C) + (X - C)'B\{d(X - C)\} \\ &= (dX)'B(X - C) + (X - C)'B dX \end{aligned}$$

From this we obtain

$$\begin{aligned} d \text{vec}(G) &= \{(X - C)'B \otimes I_n\} \text{vec}(dX') + \{I_n \otimes (X - C)'B\} \text{vec}(dX) \\ &= \{(X - C)'B \otimes I_n\} K_{mn} \text{vec}(dX) + \{I_n \otimes (X - C)'B\} \text{vec}(dX) \\ &= K_{nn} \{I_n \otimes (X - C)'B\} \text{vec}(dX) + \{I_n \otimes (X - C)'B\} \text{vec}(dX) \\ &= (I_n^2 + K_{nn}) \{I_n \otimes (X - C)'B\} \text{vec}(dX) \\ &= 2N_n \{I_n \otimes (X - C)'B\} d \text{vec}(X), \end{aligned}$$

where we have used properties of the vec operator and the commutation matrix. Consequently, we have

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(G) = 2N_n \{I_n \otimes (X - C)'B\}$$

In our next example, we show how the Jacobian matrix of the simple transformation  $z = c + Ax$  can be used to obtain the multivariate normal density function given in (1.13).

**Example 8.3.** Suppose that  $z$  is an  $m \times 1$  random vector with density function  $f_1(z)$  that is positive for all  $z \in S_1 \subseteq R^m$ . Let the  $m \times 1$  vector  $x = x(z)$  represent a one-to-one transformation of  $S_1$  onto  $S_2 \subseteq R^m$ , so that the inverse transformation  $z = z(x)$ ,  $x \in S_2$  is unique. Denote the Jacobian matrix of  $z$  at  $x$  as

$$J = \frac{\partial}{\partial x'} z(x)$$

If the partial derivatives in  $J$  exist and are continuous functions on the set  $S_2$ , then the density of  $x$  is given by

$$f_2(x) = f_1(z(x)) |J|$$

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

$$\begin{aligned} f_2(\mathbf{x}) &= \frac{1}{(2\pi)^{m/2}} \exp\left(-\frac{1}{2} \{T^{-1}(\mathbf{x} - \boldsymbol{\mu})\}' T^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) |T^{-1}| \\ &= \frac{1}{(2\pi)^{m/2} |T|} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' T'^{-1} T^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \\ &= \frac{1}{(2\pi)^{m/2} |\Omega|^{1/2}} \exp\left(-\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})' \Omega^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \end{aligned}$$

#### 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 8.1.** Let  $X$  be an  $m \times m$  matrix. Then

$$(a) \quad d\{\text{tr}(X)\} = \text{vec}(\mathbf{I}_m)' d \text{vec}(X); \quad \frac{\partial}{\partial \text{vec}(X)'} \text{tr}(X) = \text{vec}(\mathbf{I}_m)',$$

$$(b) \quad d|X| = \text{tr}(X_{\#} dX) = |X| \text{tr}(X^{-1} dX); \quad \frac{\partial}{\partial \text{vec}(X)'} |X| = \text{vec}(X'_{\#})',$$

where  $X_{\#}$  is the adjoint matrix of  $X$ .

*Proof.* Part (a) follows directly from the fact that

$$d \text{tr}(X) = \text{tr}(dX) = \text{tr}(\mathbf{I}_m dX) = \text{vec}(\mathbf{I}_m)' \text{vec}(dX) = \text{vec}(\mathbf{I}_m)' d \text{vec}(X),$$



with the third equality following from Theorem 7.15. 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 and the fact that  $X^{-1} = |X|^{-1} X_{\#}$ , we also get

$$d|X| = \{\text{vec}(X'_{\#})\}' \text{vec}(dX) = \text{tr}(X_{\#} dX) = |X| \text{tr}(X^{-1} dX). \quad \square$$

An immediate consequence of Theorem 8.1(b) is the following.

**Corollary 8.1.1.** Let  $X$  be an  $m \times m$  nonsingular matrix. Then

$$d\{\log(|X|)\} = \text{tr}(X^{-1} dX); \quad \frac{\partial}{\partial \text{vec}(X)'} \log(|X|) = \text{vec}(X^{-1})'$$

Our next result gives the differential and derivative of the inverse of a nonsingular matrix.

**Theorem 8.2.** 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

$$dX^{-1} = -X^{-1}(dX)X^{-1}$$

The expression given for the derivative now follows since

$$\begin{aligned} d \operatorname{vec}(X^{-1}) &= \operatorname{vec}(dX^{-1}) = -\operatorname{vec}(X^{-1}(dX)X^{-1}) \\ &= -(X^{-1'} \otimes X^{-1})\operatorname{vec}(dX) = -(X^{-1'} \otimes X^{-1}) d \operatorname{vec}(X) \quad \square \end{aligned}$$

A natural generalization of Theorem 8.2 is one that gives the differential and derivative of the Moore–Penrose inverse of a matrix. The following theorem gives the form of these when they exist at a matrix  $X$ .

**Theorem 8.3.** If  $X$  is an  $m \times n$  matrix and  $X^+$  is its Moore–Penrose inverse, then

$$dX^+ = (I - X^+X)(dX')X^{+'}X^+ + X^+X^{+'}(dX')(I - XX^+) - X^+(dX)X^+$$

and

$$\begin{aligned} \frac{\partial}{\partial \operatorname{vec}(X)'} \operatorname{vec}(X^+) &= \{X^{+'}X^+ \otimes (I - X^+X) + (I - XX^+) \otimes X^+X^{+'}\} \\ &\quad \cdot K_{mn} - (X^{+'} \otimes X^+) \end{aligned}$$

*Proof.* Note that

$$d(XX^+) = (dX)X^+ + X dX^+,$$

from which we get

$$XdX^+ = d(XX^+) - (dX)X^+ \tag{8.5}$$

Since  $X^+ = X^+XX^+$ , we also have

$$\begin{aligned} dX^+ &= d(X^+XX^+) = d(X^+X)X^+ + X^+X dX^+ \\ &= d(X^+X)X^+ + X^+ d(XX^+) - X^+(dX)X^+, \end{aligned} \tag{8.6}$$

where we have used (8.5) 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  $d(XX^+)$ ,

use the fact that  $XX^+$  is symmetric and idempotent to get

$$\begin{aligned} d(XX^+) &= d(XX^+XX^+) = d(XX^+)XX^+ + XX^+d(XX^+) \\ &= d(XX^+)XX^+ + (d(XX^+)XX^+)', \end{aligned} \quad (8.7)$$

since  $d(XX^+)' = d((XX^+)') = d(XX^+)$ . But

$$d(XX^+)X = dX - XX^+ dX = (I - XX^+) dX, \quad (8.8)$$

since  $X = XX^+X$  implies that

$$dX = d(XX^+X) = d(XX^+)X + XX^+ dX$$

Now substituting (8.8) in (8.7), we find that

$$\begin{aligned} dXX^+ &= (I - XX^+)(dX)X^+ + \{(I - XX^+)(dX)X^+\}' \\ &= (I - XX^+)(dX)X^+ + X^{+'}(dX')(I - XX^+) \end{aligned} \quad (8.9)$$

By using the fact that  $X^+X$  is symmetric and idempotent, we can show in a similar fashion that

$$dX^+X = X^+(dX)(I - X^+X) + (I - X^+X)(dX')X^{+'} \quad (8.10)$$

Substituting (8.9) and (8.10) into (8.6) and noting that  $(I - X^+X)X^+ = (0)$  and  $X^+(I - XX^+) = (0)$ , we get

$$dX^+ = (I - X^+X)(dX')X^{+'}X^+ + X^+X^{+'}(dX')(I - XX^+) - X^+(dX)X^+,$$

as is required. The expression given for the derivative follows since, when we take the  $\text{vec}$  of both sides of the equation above, we get

$$\begin{aligned} d \text{vec}(X^+) &= \{X^{+'}X^+ \otimes (I - X^+X)\} \text{vec}(dX') + \{(I - XX^+) \otimes X^+X^{+'}\} \text{vec}(dX') \\ &\quad - (X^{+'} \otimes X^+) \text{vec}(dX) \\ &= \{X^{+'}X^+ \otimes (I - X^+X) + (I - XX^+) \otimes X^+X^{+'}\} K_{mn} d \text{vec}(X) \\ &\quad - (X^{+'} \otimes X^+) d \text{vec}(X) \quad \square \end{aligned}$$

## 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 due to the symmetry it only contains  $m(m+1)/2$  mathematically independent variables. These variables are precisely the variables comprising the vector  $v(X)$ . If  $f(X)$  is some vector function of the matrix  $X$ , then the derivative of  $f$  will be given by the matrix

$$\frac{\partial}{\partial v(X)'} f(X)$$

We can compute derivatives of this form by utilizing the derivative

$$\frac{\partial}{\partial \text{vec}(X)'} f(X),$$

for a general nonsymmetric matrix  $X$ , along with the chain rule. Specifically, from the chain rule we have

$$\frac{\partial}{\partial v(X)'} f(X) = \left( \frac{\partial}{\partial \text{vec}(X)'} f(X) \right) \left( \frac{\partial}{\partial v(X)'} \text{vec}(X) \right)$$

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 by making use of the duplication matrix  $D_m$ . Since  $D_m v(X) = \text{vec}(X)$ , we immediately get  $D_m d v(X) = d \text{vec}(X)$ , and so

$$\frac{\partial}{\partial v(X)'} f(X) = \left( \frac{\partial}{\partial \text{vec}(X)'} f(X) \right) D_m$$

Consequently, the following results follow directly from Theorems 8.1–8.3.

**Theorem 8.4.** Let  $X$  be an  $m \times m$  symmetric matrix of variables. Then

$$(a) \quad \frac{\partial}{\partial v(X)'} |X| = \text{vec}(X'_{\#})' D_m,$$

$$(b) \quad \frac{\partial}{\partial v(X)'} \text{vec}(X^{-1}) = -(X^{-1} \otimes X^{-1}) D_m,$$

$$(c) \quad \frac{\partial}{\partial v(X)'} \text{vec}(X^+) = (\{X^+ X^+ \otimes (I - X^+ X) + (I - X X^+) \otimes X^+ X^+\}$$

$$\cdot K_{mm} - (X^+ \otimes X^+)) D_m$$

The derivatives given in (b) and (c) of Theorem 8.4 still have some redundant elements due to 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

$$\frac{\partial}{\partial v(X)'} v\{F(X)\}$$

This matrix derivative can be easily computed from the derivative

$$A = \frac{\partial}{\partial v(X)'} \text{vec}\{F(X)\}, \quad (8.11)$$

by again using the relationship  $\text{vec}(F) = D_m v(F)$ . Thus, since (8.11) implies that  $d \text{vec}(F) = A d v(X)$ , we have

$$D_m d v(F) = A d v(X),$$

or

$$D_m^+ D_m d v(F) = d v(F) = D_m^+ A d v(X),$$

since  $D_m^+ D_m = I$  by Theorem 7.36. Using this we obtain the following derivatives.

**Corollary 8.4.1.** Let  $X$  be an  $m \times m$  symmetric matrix of variables. Then

$$(a) \frac{\partial}{\partial v(X)'} v(X^{-1}) = -D_m^+(X^{-1} \otimes X^{-1})D_m,$$

$$(b) \frac{\partial}{\partial v(X)'} v(X^+) = D_m^+ (\{X^+ X^+ \otimes (I - X^+ X) + (I - X X^+) \otimes X^+ X^+\} \\ \cdot K_{mm} - (X^+ \otimes X^+)) D_m.$$

## 6. THE PERTURBATION METHOD

The perturbation method is a technique, closely related to the method utilizing 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  $\epsilon$  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  $f$  of  $X$  would then be of the form

$$f(X + \epsilon Y) = f(X) + \sum_{i=1}^{\infty} \epsilon^i g_i(X, Y),$$

where  $g_i(X, Y)$  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

$$F(X + \epsilon Y) = F(X) + \sum_{i=1}^{\infty} \epsilon^i G_i(X, Y) \quad (8.12)$$

Our goal is to determine the first few terms in the summations given above. These then can be used in an approximation of  $f(X + \epsilon Y)$  or  $F(X + \epsilon Y)$  when  $\epsilon$  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 (8.12) can be written

$$(X + \epsilon Y)^{-1} = X^{-1} + \epsilon G_1 + \epsilon^2 G_2 + \epsilon^3 G_3 + \dots$$

But we must have

$$\begin{aligned} I_m &= (X + \epsilon Y)(X + \epsilon Y)^{-1} \\ &= (X + \epsilon Y)(X^{-1} + \epsilon G_1 + \epsilon^2 G_2 + \epsilon^3 G_3 + \dots) \\ &= I_m + \epsilon(YX^{-1} + XG_1) + \epsilon^2(YG_1 + XG_2) + \epsilon^3(YG_2 + XG_3) + \dots \end{aligned}$$

If this is to hold for all  $\epsilon$ , 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

$$G_h = -X^{-1} Y G_{h-1}$$

As a result, we have

$$(X + \epsilon Y)^{-1} = X^{-1} - \epsilon X^{-1} Y X^{-1} + \epsilon^2 X^{-1} Y X^{-1} Y X^{-1} - \epsilon^3 X^{-1} Y X^{-1} Y X^{-1} Y X^{-1} + \dots,$$

or, if we return to the notation  $dX = \epsilon Y$ ,

$$(X + dX)^{-1} = X^{-1} - X^{-1} (dX) X^{-1} + X^{-1} (dX) X^{-1} (dX) X^{-1} - X^{-1} (dX) X^{-1} (dX) X^{-1} (dX) X^{-1} + \dots$$

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 8.5.** Suppose  $X = \text{diag}(x_1, \dots, x_m)$ , where  $x_1 \geq \dots \geq x_{l-1} > x_l > x_{l+1} \geq \dots \geq x_m$ , 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

$$\begin{aligned} \lambda_l(X + U) &\approx x_l + u_{ll} - \sum_{i \neq l} \frac{u_{il}^2}{(x_i - x_l)} - \sum_{i \neq l} \frac{u_{ll} u_{il}^2}{(x_i - x_l)^2} \\ &\quad + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il} u_{jl} u_{ij}}{(x_i - x_l)(x_j - x_l)}, \\ \gamma_{ll}(X + U) &\approx 1 - \frac{1}{2} \sum_{i \neq l} \frac{u_{il}^2}{(x_i - x_l)^2} - \sum_{i \neq l} \frac{u_{ll} u_{il}^2}{(x_i - x_l)^3} \\ &\quad + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il} u_{jl} u_{ij}}{(x_i - x_l)^2 (x_j - x_l)}, \end{aligned}$$

and for  $h \neq l$

$$\begin{aligned}
\gamma_{hl}(X+U) \approx & -\frac{u_{hl}}{(x_h-x_l)} - \frac{u_{ll}u_{hl}}{(x_h-x_l)^2} + \sum_{i \neq l} \frac{u_{il}u_{hi}}{(x_h-x_l)(x_i-x_l)} - \frac{u_{ll}^2 u_{hl}}{(x_h-x_l)^3} \\
& + \sum_{i \neq l} \frac{u_{ll}u_{hi}u_{il}}{(x_h-x_l)^2(x_i-x_l)} + \sum_{i \neq l} \frac{u_{hl}u_{il}^2}{(x_h-x_l)^2(x_i-x_l)} \\
& + \sum_{i \neq l} \frac{u_{ll}u_{hi}u_{il}}{(x_h-x_l)(x_i-x_l)^2} \\
& - \sum_{i \neq l} \sum_{j \neq l} \frac{u_{hi}u_{ij}u_{jl}}{(x_h-x_l)(x_i-x_l)(x_j-x_l)} \\
& + \frac{1}{2} \sum_{i \neq l} \frac{u_{hl}u_{il}^2}{(x_h-x_l)(x_i-x_l)^2},
\end{aligned}$$

where  $\gamma_{hl}(X+U)$  denotes the  $h$ th element of  $\gamma_l(X+U)$ , and the approximations above 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

$$\lambda_l = x_l + a_1 + a_2 + a_3 + \cdots, \quad (8.13)$$

$$\gamma_l = e_l + b_1 + b_2 + b_3 + \cdots, \quad (8.14)$$

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

$$Xe_l = x_l e_l, \quad (8.15)$$

$$Xb_1 + Ue_l = x_l b_1 + a_1 e_l, \quad (8.16)$$

$$Xb_2 + Ub_1 = x_l b_2 + a_1 b_1 + a_2 e_l, \quad (8.17)$$

$$Xb_3 + Ub_2 = x_l b_3 + a_1 b_2 + a_2 b_1 + a_3 e_l \quad (8.18)$$

In a similar fashion, the normalizing equation  $\gamma_l' \gamma_l = 1$  yields the identities

$$e_l' e_l = 1, \quad (8.19)$$

$$e_l' b_1 + b_1' e_l = 0, \quad (8.20)$$

$$e_l' b_2 + b_1' b_1 + b_2' e_l = 0, \quad (8.21)$$

$$e_l' b_3 + b_1' b_2 + b_2' b_1 + b_3' e_l = 0. \quad (8.22)$$



Equations (8.15) and (8.19) are trivially true, while equations (8.16) and (8.20) can be used to find  $a_1$  and  $b_1$ . Premultiplying (8.16) by  $e_1'$  and then solving for  $a_1$ , we find that

$$a_1 = e_1' U e_1 = u_{11} \quad (8.23)$$

We can then rewrite (8.16) as the system of linear equations

$$(X - x_1 I_m) b_1 = -(U - u_{11} I_m) e_1,$$

with the general solution for  $b_1$  given by

$$b_1 = -(X - x_1 I_m)^+ (U - u_{11} I_m) e_1 + c_1 e_1,$$

where  $c_1$  is an arbitrary constant. Since  $(X - x_1 I_m)^+ e_1 = 0$  and (8.20) implies that  $e_1' b_1 = 0$ , it follows that  $c_1 = 0$  and thus,

$$b_1 = -(X - x_1 I_m)^+ U e_1 \quad (8.24)$$

Next, we will use (8.17) and (8.21) to find  $a_2$  and  $b_2$ . Premultiplying (8.17) by  $e_1'$  and then solving for  $a_2$ , we find, after again using the fact that  $e_1' b_1 = 0$ , that

$$a_2 = e_1' U b_1 = -e_1' U (X - x_1 I_m)^+ U e_1 \quad (8.25)$$

Rewriting (8.17) as the system of equations in  $b_2$ ,

$$(X - x_1 I_m) b_2 = a_2 e_1 - (U - a_1 I_m) b_1,$$

which for any scalar  $c_2$  has as a solution

$$b_2 = (X - x_1 I_m)^+ \{a_2 e_1 - (U - a_1 I_m) b_1\} + c_2 e_1$$

Now since  $(X - x_1 I_m)^+ e_1 = 0$  and (8.21) implies that  $e_1' b_2 = -\frac{1}{2} b_1' b_1$ , we find that

$$c_2 = -\frac{1}{2} b_1' b_1 = -\frac{1}{2} e_1' U \{(X - x_1 I_m)^+\}^2 U e_1 = -\frac{1}{2} \sum_{i \neq 1} \frac{u_{i1}^2}{(x_i - x_1)^2},$$

and so with this value for  $c_2$ , the solution for  $b_2$  is given by

$$b_2 = (X - x_1 I_m)^+ (U - u_{11} I_m) (X - x_1 I_m)^+ U e_1 + c_2 e_1 \quad (8.26)$$

To find  $a_3$ , premultiply (8.18) by  $e'_l$  and solve for  $a_3$ , after using  $e'_l b_1 = 0$ , to get

$$\begin{aligned} a_3 &= e'_l (U - a_1 I_m) b_2 \\ &= e'_l (U - u_{ll} I_m) \{ (X - x_l I_m)^+ (U - u_{ll} I_m) (X - x_l I_m)^+ U e_l + c_2 e_l \} \\ &= e'_l U (X - x_l I_m)^+ (U - u_{ll} I_m) (X - x_l I_m)^+ U e_l \end{aligned} \quad (8.27)$$

Equation (8.18) can be expressed as

$$(X - x_l I_m) b_3 = a_3 e_l + a_2 b_1 - (U - a_1 I_m) b_2,$$

so that the solution for  $b_3$  will be given by

$$\begin{aligned} b_3 &= (X - x_l I_m)^+ \{ a_3 e_l + a_2 b_1 - (U - a_1 I_m) b_2 \} + c_3 e_l \\ &= (X - x_l I_m)^+ \{ a_2 b_1 - (U - a_1 I_m) b_2 \} + c_3 e_l, \end{aligned} \quad (8.28)$$

where  $c_3$  is an arbitrary constant. By premultiplying this equation by  $e'_l$  and using  $e'_l b_3 = -b'_1 b_2$ , which follows from (8.22), we find that

$$\begin{aligned} c_3 &= -b'_1 b_2 = e'_l U \{ (X - x_l I_m)^+ \}^2 (U - u_{ll} I_m) (X - x_l I_m)^+ U e_l \\ &= - \sum_{i \neq l} \frac{u_{ll} u_{il}^2}{(x_i - x_l)^3} + \sum_{i \neq l} \sum_{j \neq l} \frac{u_{il} u_{jl} u_{ij}}{(x_i - x_l)^2 (x_j - x_l)} \end{aligned}$$

The results now follow by substituting (8.23), (8.25), and (8.27) in (8.13) and (8.24), (8.26), and (8.28) in (8.14).  $\square$

Theorem 8.5 can be used 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  $q_l$ , which is the  $l$ th column of  $Q$ . As in Theorem 8.5, 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 8.5, 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

$$\begin{aligned} \lambda_l(Z + W) &\approx x_l + q_l'Wq_l, \\ \gamma_l(Z + W) &\approx Q\{e_l - (X - x_lI_m)^+(Q'WQ)e_l\} \\ &= q_l - (Z - x_lI_m)^+Wq_l \end{aligned}$$

The following is an immediate consequence of the first-order Taylor expansion formulas given above.

**Theorem 8.6.** Let  $\lambda_l(Z)$  be the eigenvalue 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

$$\begin{aligned} d\lambda_l &= \gamma_l'(dZ)\gamma_l, & \frac{\partial}{\partial v(Z)'} \lambda_l(Z) &= (\gamma_l' \otimes \gamma_l)D_m, \\ d\gamma_l &= -(Z - \lambda_lI_m)^+(dZ)\gamma_l, & \frac{\partial}{\partial v(Z)'} \gamma_l(Z) &= -\{\gamma_l' \otimes (Z - \lambda_lI_m)^+\}D_m \end{aligned}$$

The expansions given in and immediately following Theorem 8.5 do not hold when the eigenvalue  $x_l$  is not distinct. Suppose, for instance, that again  $x_1 \geq \dots \geq x_m$ , but now  $x_l = x_{l+1} = \dots = x_{l+r-1}$ , 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  $\Phi_l$  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 below. The proof, which is similar to that of Theorem 8.5, is left to the reader.

**Theorem 8.7.** Let  $Z$  be an  $m \times m$  symmetric matrix with eigenvalues  $x_1 \geq \dots \geq x_{l-1} > x_l = x_{l+1} = \dots = x_{l+r-1} > x_{l+r} \geq \dots \geq x_m$ , so that  $x_l$  is an eigenvalue with multiplicity  $r$ . Suppose that  $W$  is an  $m \times m$  symmetric matrix and let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_m$  be the eigenvalues of  $Z+W$ , while  $\bar{\lambda}_{l,l+r-1} = r^{-1}(\lambda_l + \dots + \lambda_{l+r-1})$ . 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  $\Phi_l$ . 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

$$a_1 = \frac{1}{r} \operatorname{tr}(WP_l),$$

$$a_2 = -\frac{1}{r} \operatorname{tr}(WYWP_l),$$

$$a_3 = \frac{1}{r} \{\operatorname{tr}(YWP_lW) - \operatorname{tr}(Y^2WP_lWP_lW)\},$$

$$B_1 = -YWP_l - P_lWY,$$

$$B_2 = YWP_lWY + YWYWP_l - Y^2WP_lWP_l + P_lWYWY - P_lWP_lWY^2 - P_lWY^2WP_l,$$

$$\begin{aligned}B_3 = & Y^2WP_lWYWP_l + P_lWYWP_lWY^2 + Y^2WP_lWP_lWY + YWP_lWP_lWY^2 \\ & + Y^2WYWP_lWP_l + P_lWP_lWYWY^2 + YWY^2WP_lWP_l + P_lWP_lWY^2WY \\ & - Y^3WP_lWP_lWP_l - P_lWP_lWP_lWY^3 - YWYWP_lWY - YWP_lWYWY \\ & - YWYWYWP_l - P_lWYWYWY + YWP_lWY^2WP_l + P_lWY^2WP_lWY \\ & + P_lWY^2WYWP_l + P_lWYWY^2WP_l - P_lWY^3WP_lWP_l \\ & - P_lWP_lWY^3WP_l\end{aligned}$$

## 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 the reader is referred to Khuri (1993) or Magnus and Neudecker (1988). Our first result gives a necessary condition for a function  $f$  to have a local maximum at  $a$ .

**Theorem 8.8.** Suppose the function  $f(x)$  is defined for all  $n \times 1$  vectors  $x \in S$ , where  $S$  is some subset of  $R^n$ . Let  $a$  be an interior point of  $S$ ; that is, there exists a  $\delta > 0$  such that  $a + u \in S$  for all  $u'u < \delta$ . If  $f$  has a local maximum at  $a$  and  $f$  is differentiable at  $a$ , then

$$\frac{\partial}{\partial a'} f(a) = 0' \quad (8.29)$$

Any point  $a$  satisfying (8.29) is called a stationary point of  $f$ . While Theorem 8.8 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. Our next result 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 8.9.** Suppose the function  $f(x)$  is defined for all  $n \times 1$  vectors  $x \in S$ , where  $S$  is some subset of  $R^n$ . Suppose also that  $f$  is twice differentiable at the interior point  $a$  of  $S$ . If  $a$  is a stationary point of  $f$  and  $H_f$  is the Hessian matrix of  $f$  at  $a$ , then

- (a)  $f$  has a local minimum at  $a$  if  $H_f$  is positive definite,
- (b)  $f$  has a local maximum at  $a$  if  $H_f$  is negative definite,
- (c)  $f$  has a saddle point at  $a$  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  $a$  if  $H_f$  is singular.

**Example 8.4.** On several occasions, we have discussed the problem of finding a least squares solution  $\hat{\beta}$  to the inconsistent system of equations

$$y = X\beta,$$

where  $y$  is an  $N \times 1$  vector of constants,  $X$  is an  $N \times (k+1)$  matrix of constants, and  $\beta$  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, while 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 be used to obtain

a solution. We will assume that  $\text{rank}(X) = k \times 1$ ; that is, the matrix  $X$  has full column rank. Recall that a least squares solution  $\hat{\beta}$  is any vector which minimizes the sum of squared errors given by

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

The differential of  $f(\hat{\beta})$  is

$$\begin{aligned} df &= d\{(y - X\hat{\beta})'\}(y - X\hat{\beta}) + (y - X\hat{\beta})' d(y - X\hat{\beta}) \\ &= -(d\hat{\beta})'X'(y - X\hat{\beta}) - (y - X\hat{\beta})'X d\hat{\beta} = -2(y - X\hat{\beta})'X d\hat{\beta}, \end{aligned}$$

so that

$$\frac{\partial}{\partial \hat{\beta}'} f(\hat{\beta}) = -2(y - X\hat{\beta})'X$$

Thus, upon setting this first derivative equal to  $0'$  and rearranging, we find that the stationary values are given by the solutions  $\hat{\beta}$  to the system of equations

$$X'X\hat{\beta} = X'y \quad (8.30)$$

Since  $X$  has full column rank,  $X'X$  is nonsingular, and so the unique solution to (8.30) is

$$\hat{\beta} = (X'X)^{-1}X'y \quad (8.31)$$

In order 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

$$\begin{aligned} d^2f &= d(df) = -d\{2(y - X\hat{\beta})'X d\hat{\beta}\} = -2\{d(y - X\hat{\beta})\}'X d\hat{\beta} \\ &= 2(d\hat{\beta})'X'X d\hat{\beta}, \end{aligned}$$

so that

$$H_f = 2X'X$$

Since this matrix is positive definite, it follows from Theorem 8.9 that the solution given in (8.31) minimizes  $f(\hat{\beta})$ .

**Example 8.5.** 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  $x_1, \dots, x_n$  from a population having density function  $f(x; \theta)$ , where  $\theta$  is a vector of parameters, then the likelihood function of  $\theta$  is defined to be the joint density function of  $x_1, \dots, x_n$  viewed as a function of  $\theta$ ; that is, this likelihood function is given by

$$L(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

The method of maximum likelihood estimates  $\theta$  by the vector  $\hat{\theta}$ , which maximizes  $L(\theta)$ . In this example, we will use this method to obtain estimates of  $\mu$  and  $\Omega$  when our sample is coming from the normal distribution,  $N_m(\mu, \Omega)$ . Thus,  $\mu$  is an  $m \times 1$  vector,  $\Omega$  is an  $m \times m$  positive definite matrix, and the required density function,  $f(x; \mu, \Omega)$  is given in (1.13). In deriving the estimates  $\hat{\mu}$  and  $\hat{\Omega}$ , we will find it a little bit easier to maximize the function  $\log(L(\mu, \Omega))$ , which is, of course, maximized at the same solution as  $L(\mu, \Omega)$ . After omitting terms from  $\log(L(\mu, \Omega))$  that do not involve  $\mu$  or  $\Omega$ , we find that we must maximize the function

$$g(\mu, \Omega) = -\frac{1}{2} n \log|\Omega| - \frac{1}{2} \text{tr}(\Omega^{-1} U),$$

where

$$U = \sum_{i=1}^n (x_i - \mu)(x_i - \mu)'$$

The first differential of  $g$  is given by

$$\begin{aligned} dg &= -\frac{1}{2} n d(\log|\Omega|) - \frac{1}{2} \text{tr}\{(d\Omega^{-1})U\} - \frac{1}{2} \text{tr}(\Omega^{-1} dU) \\ &= -\frac{1}{2} n \text{tr}(\Omega^{-1} d\Omega) + \frac{1}{2} \text{tr}\{\Omega^{-1}(d\Omega)\Omega^{-1}U\} \\ &\quad + \frac{1}{2} \text{tr}\left(\Omega^{-1} \left\{ (d\mu) \sum_{i=1}^n (x_i - \mu)' + \sum_{i=1}^n (x_i - \mu) d\mu' \right\}\right) \\ &= \frac{1}{2} \text{tr}\{(d\Omega)\Omega^{-1}(U - n\Omega)\Omega^{-1}\} \\ &\quad + \frac{1}{2} \text{tr}(\Omega^{-1} \{n(d\mu)(\bar{x} - \mu)' + n(\bar{x} - \mu) d\mu'\}) \end{aligned}$$

$$\begin{aligned}
&= \frac{1}{2} \operatorname{tr}\{(\mathrm{d}\Omega)\Omega^{-1}(U - n\Omega)\Omega^{-1}\} + n(\bar{x} - \mu)'\Omega^{-1} \mathrm{d}\mu \\
&= \frac{1}{2} \operatorname{vec}(\mathrm{d}\Omega)'(\Omega^{-1} \otimes \Omega^{-1})\operatorname{vec}(U - n\Omega) + n(\bar{x} - \mu)'\Omega^{-1} \mathrm{d}\mu,
\end{aligned}$$

where the second equality used Corollary 8.1.1 and Theorem 8.2, and the fifth used Theorem 7.17. Since  $\Omega$  is symmetric,  $\operatorname{vec}(\mathrm{d}\Omega) = \mathrm{d}\operatorname{vec}(\Omega) = D_m \mathrm{d}v(\Omega)$ , and so the differential may be reexpressed as

$$\mathrm{d}g = \frac{1}{2} \{\mathrm{d}v(\Omega)\}' D_m'(\Omega^{-1} \otimes \Omega^{-1})\operatorname{vec}(U - n\Omega) + n(\bar{x} - \mu)'\Omega^{-1} \mathrm{d}\mu, \quad (8.32)$$

and thus,

$$\frac{\partial}{\partial \mu'} g = n(\bar{x} - \mu)'\Omega^{-1}, \quad \frac{\partial}{\partial v(\Omega)'} g = \frac{1}{2} \{\operatorname{vec}(U - n\Omega)\}'(\Omega^{-1} \otimes \Omega^{-1})D_m$$

Upon equating these first derivatives to null vectors, we obtain the equations

$$\begin{aligned}
n\Omega^{-1}(\bar{x} - \mu) &= \mathbf{0}, \\
D_m'(\Omega^{-1} \otimes \Omega^{-1})\operatorname{vec}(U - n\Omega) &= \mathbf{0}
\end{aligned}$$

From the first of these two equations, we obtain the solution for  $\mu$ ,  $\hat{\mu} = \bar{x}$ , while the second can be rewritten as

$$D_m'(\Omega^{-1} \otimes \Omega^{-1})D_m v(U - n\Omega) = \mathbf{0},$$

since the symmetry of  $(U - n\Omega)$  implies that  $\operatorname{vec}(U - n\Omega) = D_m v(U - n\Omega)$ . Premultiplying this equation by  $D_m^+(\Omega \otimes \Omega)D_m'$  and using Theorem 7.38, we find that

$$v(U - n\Omega) = \mathbf{0}$$

Since  $(U - n\Omega)$  is symmetric this implies that  $(U - n\Omega) = (0)$ , and so the solution for  $\Omega$  is  $\hat{\Omega} = n^{-1}U$ . All that remains is to show that the solution  $(\hat{\mu}, \hat{\Omega})$  yields a maximum. By differentiating (8.32), we find that



$$\begin{aligned} d^2g &= \frac{1}{2} \{d v(\Omega)\}' D'_m \{d(\Omega^{-1} \otimes \Omega^{-1})\} \text{vec}(U - n\Omega) \\ &\quad + \frac{1}{2} \{d v(\Omega)\}' D'_m (\Omega^{-1} \otimes \Omega^{-1}) \text{vec}(dU - n d\Omega) \\ &\quad - n(d\mu)' \Omega^{-1} d\mu + n(\bar{x} - \mu)' (d\Omega^{-1}) d\mu \end{aligned}$$

Evaluating this at  $\mu = \bar{x}$  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

$$dU = n(d\mu)(\bar{x} - \mu)' + n(\bar{x} - \mu) d\mu'$$

also vanishes when evaluated at  $\mu = \bar{x}$ . Thus, at  $\mu = \bar{x}$  and  $\Omega = n^{-1}U$ ,

$$\begin{aligned} d^2g &= -\frac{n}{2} \{d v(\Omega)\}' D'_m (\Omega^{-1} \otimes \Omega^{-1}) D_m d v(\Omega) - n(d\mu)' \Omega^{-1} d\mu, \\ &= [d\mu' \quad \{d v(\Omega)\}'] H_g \begin{bmatrix} d\mu \\ d v(\Omega) \end{bmatrix}, \end{aligned}$$

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 since  $\Omega^{-1}$  and  $D'_m (\Omega^{-1} \otimes \Omega^{-1}) D_m$  are positive definite matrices. This then establishes that the solution  $(\hat{\mu}, \hat{\Omega}) = (\bar{x}, n^{-1}U)$  yields a maximum.

## 8. CONVEX AND CONCAVE FUNCTIONS

In Section 2.10, 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 8.1.** 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

$$f(cx_1 + (1 - c)x_2) \leq cf(x_1) + (1 - c)f(x_2)$$

for all  $x_1 \in S$ ,  $x_2 \in S$ , and  $0 \leq c \leq 1$ . If  $-f(x)$  is a convex function, then  $f(x)$  is said to be a concave function.

If  $f(x)$  is a convex function, then it is easily verified that the set defined by

$$T = \{z = (x', y) : x \in S, y \geq f(x)\}$$

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 2.27, 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

$$f(x) \geq f(a) + t(x - a) \quad (8.33)$$

The generalization of this result to arbitrary  $m$  is given below.

**Theorem 8.10.** 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  $a \in S$ , there exists an  $m \times 1$  vector  $t$  such that

$$f(x) \geq f(a) + t'(x - a) \quad (8.34)$$

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 2.27 that there exists an  $(m + 1) \times 1$  vector  $b = (b'_1, b_{m+1})' \neq 0$  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 since if it were, we would be able to 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 (8.34) 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 (8.34) will be given by the tangent hyperplane to  $f(\mathbf{x})$  at  $\mathbf{x} = \mathbf{a}$ .

**Theorem 8.11.** 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  $\mathbf{a} \in S$ , then

$$f(\mathbf{x}) \geq f(\mathbf{a}) + \left( \frac{\partial}{\partial \mathbf{a}'} f(\mathbf{a}) \right) (\mathbf{x} - \mathbf{a})$$

for all  $\mathbf{x} \in S$ .

*Proof.* Suppose that  $\mathbf{x} \in S$  and  $\mathbf{a} \in S$ , and let  $\mathbf{y} = \mathbf{a} - \mathbf{x}$  so that  $\mathbf{a} = \mathbf{x} + \mathbf{y}$ . Since  $S$  is convex, the point

$$c\mathbf{a} + (1 - c)\mathbf{x} = c(\mathbf{x} + \mathbf{y}) + (1 - c)\mathbf{x} = \mathbf{x} + c\mathbf{y}$$

is in  $S$  for  $0 \leq c \leq 1$ . Thus, due to the convexity of  $f$ , we have

$$f(\mathbf{x} + c\mathbf{y}) \leq cf(\mathbf{x} + \mathbf{y}) + (1 - c)f(\mathbf{x}) = f(\mathbf{x}) + c\{f(\mathbf{x} + \mathbf{y}) - f(\mathbf{x})\},$$

or, equivalently,

$$f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{x}) + c^{-1}\{f(\mathbf{x} + c\mathbf{y}) - f(\mathbf{x})\} \quad (8.35)$$

Now since  $f$  is differentiable, we also have the Taylor formula

$$f(\mathbf{x} + c\mathbf{y}) = f(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) \right) c\mathbf{y} + r_1(c\mathbf{y}, \mathbf{x}), \quad (8.36)$$

where the remainder satisfies  $\lim c^{-1}r_1(c\mathbf{y}, \mathbf{x}) = 0$  as  $c \rightarrow 0$ . Using (8.36) in (8.35), we get

$$f(\mathbf{x} + \mathbf{y}) \geq f(\mathbf{x}) + \left( \frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) \right) \mathbf{y} + c^{-1}r_1(c\mathbf{y}, \mathbf{x}),$$

and so the result follows by letting  $c \rightarrow 0$ . □

The previous theorem can easily be used 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 8.12.** 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  $a \in S$  is a stationary point of  $f$ , then  $f$  has an absolute minimum at  $a$ .

*Proof.* If  $a$  is a stationary point of  $f$ , then

$$\frac{\partial}{\partial a'} f(a) = \mathbf{0}'$$

Using this in the inequality of Theorem 8.11, we get  $f(x) \geq f(a)$  for all  $x \in S$ , and so the result follows.  $\square$

The inequality given in (8.34) can be used to prove a very useful inequality involving the moments of a random vector  $x$ . This inequality is known as Jensen's inequality. But before we can prove this result, we will need the following.

**Theorem 8.13.** Suppose that  $S$  is a convex subset of  $R^m$  and  $y$  is an  $m \times 1$  random vector with finite first moments. If  $P(y \in S) = 1$ , then  $E(y) \in S$ .

*Proof.* We will prove the result by induction. Clearly, the result holds if  $m = 1$ , since in this case  $S$  is an interval, and it is easily demonstrated 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  $S_* = \{x: x = u - E(y), u \in S\}$  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 2.27 that there exists an  $m \times 1$  vector  $a \neq \mathbf{0}$  such that  $a'x \geq 0$  for all  $x \in S_*$ . Consequently, since  $P(y \in S) = P(w \in S_*) = 1$ , where the random vector  $w = y - E(y)$ , we have  $a'w \geq 0$  with probability 1, yet  $E(a'w) = 0$ . This is possible only if  $a'w = 0$ , in which case  $w$  is on the hyperplane defined by  $\{x: a'x = 0\}$ , with probability one. But since  $P(w \in S_*) = 1$  as well, we must have  $P(w \in S_0) = 1$ , where  $S_0 = S_* \cap \{x: a'x = 0\}$ . Now it follows from Theorem 2.23 that  $S_0$  is a convex set, and it is contained within an  $(m - 1)$ -dimensional vector space since  $\{x: a'x = 0\}$  is an  $(m - 1)$ -dimensional vector space. Thus, since our result holds for  $m - 1$ -dimensional spaces, we must have  $E(w) = \mathbf{0} \in S_0$ . This leads to the contradiction  $\mathbf{0} \in S_*$ , since  $S_0 \subseteq S_*$ , and so the proof is complete.  $\square$

We now prove Jensen's inequality.

**Theorem 8.14.** 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  $y$  is an  $m \times 1$  random vector with finite first moments and satisfying  $P(y \in S) = 1$ , then

$$E(f(\mathbf{y})) \geq f(E(\mathbf{y}))$$

*Proof.* The previous theorem guarantees that  $E(\mathbf{y}) \in S$ . We first prove the result for  $m = 1$ . If  $E(\mathbf{y})$  is an interior point of  $S$ , the result follows by taking the expected value of both sides of (8.33) when  $x = \mathbf{y}$  and  $a = E(\mathbf{y})$ . Since when  $m = 1$ ,  $S$  is an interval,  $E(\mathbf{y})$  can be a boundary point of  $S$  only if  $S$  is closed and  $P(\mathbf{y} = c) = 1$ , where  $c$  is an endpoint of the interval. In this case, the result is trivial since 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 (8.34) with  $x = \mathbf{y}$  and  $a = E(\mathbf{y})$ . If  $E(\mathbf{y})$  is a boundary point of  $S$ , then we know from the supporting hyperplane theorem that there exists an  $m \times 1$  unit vector  $\mathbf{b}$  such that  $w = \mathbf{b}'\mathbf{y} \geq \mathbf{b}'E(\mathbf{y}) = \mu$  with probability one. But since 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

$$g(\mathbf{u}_1) = f\left(P\begin{pmatrix} \mathbf{u}_1 \\ \mu \end{pmatrix}\right) = f(\mathbf{y}),$$

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, since  $\mathbf{u}_1$  is  $(m - 1) \times 1$ , our result applies to  $g(\mathbf{u}_1)$ . Thus, we have

$$E(f(\mathbf{y})) = E(g(\mathbf{u}_1)) \geq g(E(\mathbf{u}_1)) = f\left(P\begin{pmatrix} E(\mathbf{u}_1) \\ \mu \end{pmatrix}\right) = f(E(\mathbf{y})) \quad \square$$

## 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  $\mathbf{x} \in S$ , while 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, there exist functions  $g_1, \dots, g_m$  such that

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

where  $g(\mathbf{x})$  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

$$\begin{aligned} \frac{\partial}{\partial \mathbf{x}'} L(\mathbf{x}, \boldsymbol{\lambda}) &= \frac{\partial}{\partial \mathbf{x}'} f(\mathbf{x}) - \boldsymbol{\lambda}' \left( \frac{\partial}{\partial \mathbf{x}'} \mathbf{g}(\mathbf{x}) \right) = \mathbf{0}', & (8.37) \\ \frac{\partial}{\partial \boldsymbol{\lambda}'} L(\mathbf{x}, \boldsymbol{\lambda}) &= -\mathbf{g}(\mathbf{x})' = \mathbf{0}' \end{aligned}$$

The second equation above is simply the equality constraints

$$\mathbf{g}(\mathbf{x}) = \mathbf{0} \quad (8.38)$$

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 (8.37) and (8.38). 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 (1988).

**Theorem 8.15.** 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

$$\mathbf{x}'\mathbf{A}\mathbf{x} < 0$$

for all  $\mathbf{x} \neq \mathbf{0}$  for which  $\mathbf{B}\mathbf{x} = \mathbf{0}$ .

A similar result holds for a local minimum with the inequality  $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$  replacing  $\mathbf{x}'\mathbf{A}\mathbf{x} < 0$ . Our next result provides a method for determining whether  $\mathbf{x}'\mathbf{A}\mathbf{x} < 0$  or  $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$  holds for all  $\mathbf{x} \neq \mathbf{0}$  satisfying  $\mathbf{B}\mathbf{x} = \mathbf{0}$ . Again, a proof can be found in Magnus and Neudecker (1988).

**Theorem 8.16.** Let  $\mathbf{A}$  be an  $n \times n$  symmetric matrix and  $\mathbf{B}$  be an  $m \times n$  matrix. For  $r = 1, \dots, n$ , let  $\mathbf{A}_{rr}$  be the  $r \times r$  matrix obtained by deleting the last  $n - r$  rows and columns of  $\mathbf{A}$ , and let  $\mathbf{B}_r$  be the  $m \times r$  matrix obtained by deleting the last  $n - r$  columns of  $\mathbf{B}$ . For  $r = 1, \dots, n$ , define the  $(m+r) \times (m+r)$  matrix  $\Delta_r$  as

$$\Delta_r = \begin{bmatrix} (0) & \mathbf{B}_r \\ \mathbf{B}'_r & \mathbf{A}_{rr} \end{bmatrix}$$

Then, if  $\mathbf{B}_m$  is nonsingular,  $\mathbf{x}'\mathbf{A}\mathbf{x} > 0$  holds for all  $\mathbf{x} \neq \mathbf{0}$  satisfying  $\mathbf{B}\mathbf{x} = \mathbf{0}$  if and only if

$$(-1)^m |\Delta_r| > 0$$

for  $r = m + 1, \dots, n$ , and  $\mathbf{x}'\mathbf{A}\mathbf{x} < 0$  holds for all  $\mathbf{x} \neq \mathbf{0}$  satisfying  $\mathbf{B}\mathbf{x} = \mathbf{0}$  if and only if

$$(-1)^r |\Delta_r| > 0,$$

for  $r = m + 1, \dots, n$ .

**Example 8.6.** 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

$$x_1^2 + x_2^2 = 1, \tag{8.39}$$

$$x_3 - x_1 - x_2 = 1 \tag{8.40}$$

Setting the first derivative of the Lagrange function

$$L(x, \lambda) = x_1 + x_2 + x_3 - \lambda_1(x_1^2 + x_2^2 - 1) - \lambda_2(x_3 - x_1 - x_2 - 1),$$

with respect to  $x$ , equal to  $0'$ , 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 in (8.39), 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 Theorems 8.15 and 8.16. Thus, since  $m = 2$  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  $L(x, \lambda)$ , 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  $\Delta_r$  matrices.

**Example 8.7.** Let  $A$  be an  $m \times m$  symmetric matrix and  $x$  be an  $m \times 1$  vector. We saw in Section 3.6 that

$$\frac{x'Ax}{x'x} \tag{8.41}$$

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 since  $z = (x'x)^{-1/2}x$  is a unit vector, it follows that maximizing or minimizing (8.41) over all  $x \neq 0$  is equivalent to maximizing or minimizing the function

$$f(z) = z'Az,$$

subject to the constraint

$$z'z = 1 \tag{8.42}$$

Thus, the Lagrange function is

$$L(z, \lambda) = z'Az - \lambda(z'z - 1)$$

Setting its first derivative, with respect to  $z$ , equal to  $0'$ , we obtain the equation

$$2Az - 2\lambda z = \mathbf{0},$$

or, equivalently,

$$Az = \lambda z, \quad (8.43)$$

which is the eigenvalue–eigenvector equation. Thus, the Lagrange multiplier  $\lambda$  is an eigenvalue of  $A$ . Further, premultiplying (8.43) by  $z'$  and using (8.42), we find that

$$\lambda = z'Az;$$

that is, if  $z$  is a stationary point of  $L(z, \lambda)$ , then  $z'Az$  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  $\lambda_m(A)$ , and this is attained at any unit eigenvector associated with  $\lambda_m(A)$ .

In our final example, we obtain the best quadratic unbiased estimator of  $\sigma^2$  in the ordinary least squares regression model.

**Example 8.8.** Consider the multiple regression model  $y = X\beta + \epsilon$ , where  $\epsilon \sim N_N(\mathbf{0}, \sigma^2 I)$ . A quadratic estimator of  $\sigma^2$  is any estimator,  $\hat{\sigma}^2$  that takes the form  $\hat{\sigma}^2 = y'Ay$ , 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 since  $E(\epsilon) = \mathbf{0}$  and  $E(\epsilon\epsilon') = \sigma^2 I$ , we have

$$\begin{aligned} E(y'Ay) &= E\{(X\beta + \epsilon)'A(X\beta + \epsilon)\} = E\{\beta'X'AX\beta + 2\beta'X'A\epsilon + \epsilon'A\epsilon\} \\ &= \beta'X'AX\beta + \text{tr}\{AE(\epsilon\epsilon')\} = \beta'X'AX\beta + \sigma^2 \text{tr}(A), \end{aligned}$$

and so  $\hat{\sigma}^2 = y'Ay$  is unbiased regardless of the value of  $\beta$  only if

$$X'AX = (0) \quad (8.44)$$

and

$$\text{tr}(A) = 1 \quad (8.45)$$

Using the fact that the components of  $\epsilon$  are independently distributed and the first four moments of each component are 0, 1, 0, 3, it is easily verified that

$$\text{var}(y'Ay) = 2\sigma^4 \text{tr}(A^2) + 4\sigma^2 \beta'X'A^2X\beta$$

Thus, the required Lagrange function is

$$L(A, \lambda, \Lambda) = 2\sigma^4 \operatorname{tr}(A^2) + 4\sigma^2 \boldsymbol{\beta}' X' A^2 X \boldsymbol{\beta} - \operatorname{tr}(\Lambda X' A X) - \lambda \{ \operatorname{tr}(A) - 1 \},$$

where the Lagrange multipliers are given by  $\lambda$  and the components of the matrix  $\Lambda$ , which is symmetric since  $X' A X$  is symmetric. Differentiation with respect to  $A$  yields

$$\begin{aligned} dL &= 2\sigma^4 \operatorname{tr}\{(dA)A + A dA\} + 4\sigma^2 \boldsymbol{\beta}' X' \{(dA)A + A dA\} X \boldsymbol{\beta} \\ &\quad - \operatorname{tr}\{\Lambda X' (dA) X\} - \lambda \operatorname{tr}(dA) \\ &= \operatorname{tr}\{4\sigma^4 A + 4\sigma^2 (A X \boldsymbol{\beta} \boldsymbol{\beta}' X' + X \boldsymbol{\beta} \boldsymbol{\beta}' X' A) - X \Lambda X' - \lambda I_N\} dA \end{aligned}$$

Thus, we must use

$$4\sigma^4 A + 4\sigma^2 (A X \boldsymbol{\beta} \boldsymbol{\beta}' X' + X \boldsymbol{\beta} \boldsymbol{\beta}' X' A) - X \Lambda X' - \lambda I_N = (0) \quad (8.46)$$

along with (8.44) and (8.45) to solve for  $A$ . Premultiplying and postmultiplying (8.46) by  $XX^+$  and using (8.44) and the fact that  $X^+ = (X'X)^+ X'$ , we find that

$$X \Lambda X' = -\lambda X X^+$$

Substituting this back into (8.46), we get

$$A = \frac{1}{4} \sigma^{-4} \lambda (I_N - X X^+) - \sigma^{-2} H, \quad (8.47)$$

where  $H = A \boldsymbol{\gamma} \boldsymbol{\gamma}' + \boldsymbol{\gamma} \boldsymbol{\gamma}' A$  and  $\boldsymbol{\gamma} = X \boldsymbol{\beta}$ . Putting (8.47) back into (8.46) and simplifying, we obtain

$$H = -\sigma^{-2} (H \boldsymbol{\gamma} \boldsymbol{\gamma}' + \boldsymbol{\gamma} \boldsymbol{\gamma}' H) \quad (8.48)$$

By postmultiplying (8.48) by  $\boldsymbol{\gamma}$ , we find that  $\boldsymbol{\gamma}$  must be an eigenvector of  $H$ , and in light of equation (8.48), this can be true only if  $H$  is of the form  $H = c \boldsymbol{\gamma} \boldsymbol{\gamma}'$  for some scalar  $c$ . Further, when we put  $H = c \boldsymbol{\gamma} \boldsymbol{\gamma}'$  in (8.48), we find that we must have  $c = 0$ ; thus  $H = (0)$ . In addition, if we take the trace of both sides of (8.47) and use (8.45), we see that

$$\lambda = \frac{4\sigma^4}{\operatorname{tr}(I_N - X X^+)} = \frac{4\sigma^4}{N - r},$$

where  $r$  is the rank of  $X$ . Consequently, we have shown that (8.47) simplifies

to

$$A = (N - r)^{-1}(\mathbf{I}_N - \mathbf{X}\mathbf{X}'), \quad (8.49)$$

so that  $\hat{\sigma}^2 = \mathbf{y}'\mathbf{A}\mathbf{y} = \text{SSE}/(N - r)$  is the familiar residual variance estimate. We can easily demonstrate that (8.49) yields an absolute minimum by writing an arbitrary symmetric matrix satisfying (8.44) and (8.45), as  $A_* = A + B$ , where  $B$  must then satisfy  $\text{tr}(B) = 0$  and  $\mathbf{X}'\mathbf{B}\mathbf{X} = (0)$ . Then, since  $\text{tr}(AB) = 0$  and  $\mathbf{A}\mathbf{X} = (0)$ , we have

$$\begin{aligned} \text{var}(\mathbf{y}'A_*\mathbf{y}) &= 2\sigma^4 \text{tr}(A_*^2) + 4\sigma^2 \boldsymbol{\beta}'\mathbf{X}'A_*^2\mathbf{X}\boldsymbol{\beta} \\ &= 2\sigma^4 \{\text{tr}(A^2) + \text{tr}(B^2) + 2\text{tr}(AB)\} + 4\sigma^2 \boldsymbol{\beta}'\mathbf{X}' \\ &\quad \cdot (A^2 + B^2 + AB + BA)\mathbf{X}\boldsymbol{\beta} \\ &= 2\sigma^4 \{\text{tr}(A^2) + \text{tr}(B^2)\} + 4\sigma^2 \boldsymbol{\beta}'\mathbf{X}'B^2\mathbf{X}\boldsymbol{\beta} \\ &\geq 2\sigma^4 \text{tr}(A^2) = \text{var}(\mathbf{y}'\mathbf{A}\mathbf{y}) \end{aligned}$$

## PROBLEMS

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  $\log(1.1)$ .
2. Suppose the function  $f$  of the  $2 \times 1$  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  $u_2$ .

3. Suppose the  $2 \times 1$  function  $f$  of the  $3 \times 1$  vector  $\mathbf{x}$  is given by

$$f(\mathbf{x}) = \begin{bmatrix} x_1^2 + x_2^2 + x_3^2 \\ 2x_1 - x_2 - x_3 \end{bmatrix}$$

and the  $2 \times 1$  function  $g$  of the  $2 \times 1$  vector  $\mathbf{z}$  is given by

$$g(\mathbf{z}) = \begin{bmatrix} z_2/z_1 \\ z_1 z_2 \end{bmatrix}$$

Use the chain rule to compute

$$\frac{\partial}{\partial x'} y(x),$$

where  $y(x)$  is the composite function defined by  $y(x) = g(f(x))$ .

4. Let  $A$  and  $B$  be  $m \times m$  symmetric matrices of constants and  $x$  be an  $m \times 1$  vector of variables. Find the differential and first derivative of the function

$$f(x) = \frac{x'Ax}{x'Bx}$$

5. 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)  $\text{tr}(AXBX)$ .

6. Let  $X$  be an  $m \times m$  nonsingular matrix and  $A$  be an  $m \times m$  matrix of constants. Find the differential and derivative of

- (a)  $|X^2|$ ,  
 (b)  $\text{tr}(AX^{-1})$ .

7. Let  $X$  be an  $m \times n$  matrix with  $\text{rank}(X) = n$ . Show that

$$\frac{\partial}{\partial \text{vec}(X)'} |X'X| = 2|X'X|(\text{vec}\{X(X'X)^{-1}\})'$$

8. 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. 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)  $\text{vec}(AX^{-1}B)$ .

10. Show that if  $X$  is an  $m \times m$  nonsingular matrix and  $X_{\#}$  is its adjoint matrix, then

$$\frac{\partial}{\partial \text{vec}(X)'} \text{vec}(X\#) = |X| \{ \text{vec}(X^{-1}) \text{vec}(X^{-1}')' - (X^{-1}' \otimes X^{-1}) \}$$

11. Prove Corollary 8.1.1.

12. Let  $X$  be an  $m \times m$  symmetric matrix of variables. For each of the following functions, find the Jacobian matrix

$$\frac{\partial}{\partial v(X)'} \text{vec}(F)$$

(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.

13. 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

$$\frac{\partial}{\partial \tilde{v}(X)'} \tilde{v}(X^{-1}) = -2\tilde{L}_m(X^{-1} \otimes X^{-1})\tilde{L}_m'$$

14. Suppose that  $Y$  is an  $m \times m$  symmetric matrix and  $\epsilon$  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

$$B_1 = -\frac{1}{2}Y, B_2 = \frac{3}{8}Y^2, B_3 = -\frac{5}{16}Y^3, \text{ and } B_4 = \frac{35}{128}Y^4.$$

15. Let  $S$  be an  $m \times m$  sample covariance matrix, and suppose that  $\Omega$ , 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  $\Omega$ , while 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

$$\begin{aligned}
 C_1 &= A - \frac{1}{2} (\Omega D_A + D_A \Omega), \\
 C_2 &= \frac{3}{8} (D_A^2 \Omega + \Omega D_A^2) + \frac{1}{4} D_A \Omega D_A - \frac{1}{2} (A D_A + D_A A), \\
 C_3 &= \frac{3}{8} (D_A^2 A + A D_A^2) + \frac{1}{4} D_A A D_A - \frac{3}{16} (D_A^2 \Omega D_A + D_A \Omega D_A^2) \\
 &\quad - \frac{5}{16} (D_A^3 \Omega + \Omega D_A^3),
 \end{aligned}$$

where  $D_A = \text{diag}(a_{11}, \dots, a_{mm})$ .

16. Derive the results given in Theorem 8.7. First obtain expressions for  $B_1$ ,  $B_2$ , and  $B_3$  by utilizing 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  $\bar{\lambda}_{l,l+r-1} = r^{-1} \text{tr}\{(Z + W)\Phi_l\}$ .

17. 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  $\lambda_l = x_l + a_l$  and  $\gamma_l = c e_l + b_l$ , 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  $\gamma_l$ .

(a) Show that  $a_l = u_{ll} - x_l v_{ll}$ .

(b) Show that if  $c = 1$  and  $\gamma_l' \gamma_l = 1$ , then

$$b_{ll} = 0, \quad b_{il} = -\frac{u_{li} - x_l v_{li}}{x_i - x_l} \quad \text{for all } i \neq l,$$

where  $b_{il}$  is the  $i$ th component of the vector  $b_l$ .

(c) Show that if  $c = 1$  and  $\gamma_l'(I_m + V)\gamma_l = 1$ , then

$$b_{ll} = -\frac{1}{2} v_{ll}, \quad b_{il} = -\frac{u_{li} - x_l v_{li}}{x_i - x_l} \quad \text{for all } i \neq l$$

(d) Show that if  $c = x_l^{1/2}$  and  $\gamma_l' \gamma_l = \lambda_l$ , then

$$b_{ll} = \frac{u_{ll} - x_l v_{ll}}{2x_l^{1/2}}, \quad b_{il} = -\frac{x_l^{1/2}(u_{li} - x_l v_{li})}{x_i - x_l}, \quad \text{for all } i \neq l$$

18. Consider the function  $f$  of the  $2 \times 1$  vector  $x$  given by

$$f(x) = 2x_1^3 + x_2^3 - 6x_1 - 27x_2$$

- (a) Determine the stationary points of  $f$ .  
 (b) Identify each of the points in part (a) as a maximum, minimum, or saddle point.

19. 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)  $x_2^3 + 2x_1^2 + x_3^2 + 2x_1x_3 - 3x_2 - x_3$ .

20. Let  $a$  be an  $m \times 1$  vector and  $B$  be an  $m \times m$  symmetric matrix, each containing constants. Let  $x$  be an  $m \times 1$  vector of variables.

- (a) Show that the function

$$f(x) = x'Bx + a'x$$

has stationary solutions given by

$$x = -\frac{1}{2} B^+ a + (I - B^+ B)y,$$

where  $y$  is an arbitrary  $m \times 1$  vector.

- (b) Show that if  $B$  is nonsingular, then there is only one stationary solution. When will this solution yield a maximum or a minimum?

21. If the Hessian matrix  $H_f$  of a function  $f$  is singular at a stationary point  $x$ , then we must take a closer look at the behavior of this function in the neighborhood of the point  $x$  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 saddlepoint.

- (a)  $x_1^4 + x_2^4$ .  
 (b)  $x_1^2 x_2^2 - x_1^4 - x_2^4$ .  
 (c)  $x_1^3 - x_2^3$ .



## PROBLEMS

22. Suppose that we have independent random samples from each of  $k$  multivariate normal distributions with the  $i$ th distribution being  $N_m(\boldsymbol{\mu}_i, \boldsymbol{\Omega})$ . Thus, these distributions have possibly different mean vectors but identical covariance matrices. If the  $i$ th sample is denoted by  $x_{i1}, \dots, x_{in_i}$ , show that the maximum likelihood estimators of  $\boldsymbol{\mu}_i$  and  $\boldsymbol{\Omega}$  are given by

$$\hat{\boldsymbol{\mu}}_i = \bar{\boldsymbol{x}}_i = \sum_{j=1}^{n_i} \frac{\boldsymbol{x}_{ij}}{n_i}, \quad \hat{\boldsymbol{\Omega}} = \sum_{i=1}^k \sum_{j=1}^{n_i} \frac{(\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)(\boldsymbol{x}_{ij} - \bar{\boldsymbol{x}}_i)'}{n},$$

where  $n = n_1 + \dots + n_k$ .

23. Consider the multiple regression model,

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

where  $\boldsymbol{y}$  is  $N \times 1$ ,  $\boldsymbol{X}$  is  $N \times m$ ,  $\boldsymbol{\beta}$  is  $m \times 1$ , and  $\boldsymbol{\epsilon}$  is  $N \times 1$ . Suppose that  $\text{rank}(\boldsymbol{X}) = m$  and  $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 \boldsymbol{I}_N)$ , so that  $\boldsymbol{y} \sim N_N(\boldsymbol{X}\boldsymbol{\beta}, \sigma^2 \boldsymbol{I}_N)$ . Find the maximum likelihood estimates of  $\boldsymbol{\beta}$  and  $\sigma^2$ .

24. Let  $f(\boldsymbol{x})$  be a real-valued convex function defined for all  $\boldsymbol{x} \in S$ , where  $S$  is a convex subset of  $R^m$ . Show that the set  $T = \{\boldsymbol{z} = (\boldsymbol{x}', y)': \boldsymbol{x} \in S, y \geq f(\boldsymbol{x})\}$  is convex.
25. Suppose that  $f(\boldsymbol{x})$  and  $g(\boldsymbol{x})$  are convex functions both defined on the convex set  $S \subseteq R^m$ . Show that the function  $af(\boldsymbol{x}) + bg(\boldsymbol{x})$  is convex if  $a$  and  $b$  are nonnegative scalars.
26. Prove the converse of Theorem 8.11; that is, show that if  $f(\boldsymbol{x})$  is defined and differentiable on the open convex set  $S$  and

$$f(\boldsymbol{x}) \geq f(\boldsymbol{a}) + \left( \frac{\partial}{\partial \boldsymbol{a}'} f(\boldsymbol{a}) \right) (\boldsymbol{x} - \boldsymbol{a})$$

for all  $\boldsymbol{x} \in S$  and  $\boldsymbol{a} \in S$ , then  $f(\boldsymbol{x})$  is a convex function.

27. Let  $f(\boldsymbol{x})$  be a real-valued function defined for all  $\boldsymbol{x} \in S$ , where  $S$  is an open convex subset of  $R^m$ , and suppose that  $f(\boldsymbol{x})$  is a twice differentiable function on  $S$ . Show that  $f(\boldsymbol{x})$  is a convex function if and only if the Hessian matrix  $H_f$  is nonnegative definite at each  $\boldsymbol{x} \in S$ .

28. Let  $\boldsymbol{x}$  be a  $2 \times 1$  vector and consider the function  $f(\boldsymbol{x}) = x_1^\alpha x_2^{1-\alpha}$  for all  $\boldsymbol{x} \in S$ , where  $0 < \alpha < 1$  and  $S = \{\boldsymbol{x}: x_1 > 0, x_2 > 0\}$ .

- (a) Use the previous exercise to show that  $f(x)$  is a concave function.  
 (b) Show that if  $y$  is a  $2 \times 1$  random vector with finite first moments and satisfying  $P(y \in S) = 1$ , then

$$E(y_1^\alpha y_2^{1-\alpha}) \leq \{E(y_1)\}^\alpha \{E(y_2)\}^{1-\alpha}$$

if  $0 < \alpha < 1$ .

29. Let  $x$  be a  $3 \times 1$  vector and define the function

$$f(x) = x_1 + x_2 - x_3$$

Find the maximum and minimum of  $f(x)$  subject to the constraint  $x'x = 1$ .

30. Find the shortest distance from the origin to a point on the surface given by

$$x_1^2 + x_2^2 + x_3^2 + 4x_1 - 6x_3 = 2$$

31. Let  $A$  be an  $m \times m$  positive definite matrix and  $x$  be an  $m \times 1$  vector. Find the maximum and minimum of the function

$$f(x) = x'x,$$

subject to the constraint  $x'Ax = 1$ .

32. Find the maximum and minimum of the function

$$f(x) = x_1(x_2 + x_3),$$

subject to the constraints  $x_1^2 + x_2^2 = 1$  and  $x_1x_3 + x_2 = 2$ .

33. For a  $3 \times 1$  vector  $x$ , maximize the function

$$f(x) = x_1x_2x_3,$$

subject to the constraint  $x_1 + x_2 + x_3 = a$ , where  $a$  is some positive number. Use this to establish the inequality

$$(x_1x_2x_3)^{1/3} \leq \frac{1}{3} (x_1 + x_2 + x_3)$$

for all positive real numbers  $x_1$ ,  $x_2$ , and  $x_3$ . Generalize this result to  $m$

variables; that is, if  $x$  is  $m \times 1$ , show that

$$(x_1 x_2 \cdots x_m)^{1/m} \leq \frac{1}{m} (x_1 + \cdots + x_m)$$

holds for all positive real numbers  $x_1, \dots, x_m$ .

34. Let  $A$  and  $B$  be  $m \times m$  matrices, with  $A$  being nonnegative definite and  $B$  being positive definite. Following the approach of Example 8.7, use the Lagrange method to find the maximum and minimum values of

$$f(x) = \frac{x'Ax}{x'Bx},$$

over all  $x \neq 0$ .

35. Let  $a$  be an  $m \times 1$  vector and  $B$  be an  $m \times m$  positive definite matrix. Using the results of the previous exercise, show that for  $x \neq 0$ ,

$$f(x) = \frac{(a'x)^2}{x'Bx}$$

has a maximum value of

$$a'B^{-1}a$$

This result can be used to obtain the union–intersection test (see Example 3.14) of the multivariate hypothesis  $H_0: \mu = \mu_0$  against  $H_1: \mu \neq \mu_0$ , where  $\mu$  represents the  $m \times 1$  mean vector of a population and  $\mu_0$  is an  $m \times 1$  vector of constants. Let  $\bar{x}$  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

$$t = \frac{(\bar{x} - \mu_0)}{s/\sqrt{n}}$$

for testing  $H_0: \mu = \mu_0$ , then the union–intersection test can be based on  $T^2 = n(\bar{x} - \mu_0)'S^{-1}(\bar{x} - \mu_0)$ .

36. Suppose that  $x_1, \dots, x_n$  are independent and identically distributed random variables with mean  $\mu$  and variance  $\sigma^2$ . Consider a linear estimator of  $\mu$  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  $\mu$ ?
- (b) Use the method of Lagrange multipliers to show that the sample mean  $\bar{x}$  is the best linear unbiased estimator of  $\mu$ ; that is,  $\bar{x}$  has the smallest variance among all linear unbiased estimators of  $\mu$ .

37. 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

$$P(x_1 = n_1, \dots, x_k = n_k) = \frac{n!}{n_1! \cdots n_k!} p_1^{n_1} \cdots p_k^{n_k},$$

where  $n_1, \dots, n_k$  are nonnegative integers satisfying  $n_1 + \dots + n_k = n$ . Find the maximum likelihood estimate of  $\mathbf{p} = (p_1, \dots, p_k)'$ .

38. Suppose that the  $m \times m$  positive definite covariance matrix  $\Omega$  is partitioned in the form

$$\Omega = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega'_{12} & \Omega_{22} \end{bmatrix},$$

where  $\Omega_{11}$  is  $m_1 \times m_1$ ,  $\Omega_{22}$  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  $\Omega$  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  $f(\mathbf{x})$ , 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?

39. Consider the function,  $f(P) = \text{tr}(PXP'D)$ , where  $P$  is an  $m \times m$  orthogonal matrix, and both  $X$  and  $D$  are  $m \times m$  positive definite matrices. Further, suppose that  $D$  is diagonal with distinct, descending, positive diagonal elements; that is,  $D = \text{diag}(d_1, \dots, d_m)$  with  $d_1 > \dots > d_m > 0$ .

(a) By working with the Lagrange function,

$$L(P, \Lambda) = \text{tr}(PXP'D) + \text{tr}\{\Lambda(PP' - I_m)\},$$

where  $\Lambda$  is a symmetric matrix of Lagrange multipliers, show that the stationary points of  $f(P)$  occur when  $PXP'$  is diagonal.

(b) Use part (a) to show that

$$\max_{P: PP'=1} f(P) = \sum_{i=1}^m d_i \lambda_i(X),$$

and

$$\min_{P: PP'=1} f(P) = \sum_{i=1}^m d_{m+1-i} \lambda_i(X),$$

where  $\lambda_1(X) \geq \dots \geq \lambda_m(X) > 0$  are the eigenvalues of  $X$ .

## CHAPTER NINE

# Some Special Topics Related to Quadratic Forms

### 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, while  $A$  is a matrix of constants. The most common situation is one in which  $\mathbf{x}$  has as its distribution, or 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.

### 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 the next section 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 basic results regarding idempotent matrices.

**Theorem 9.1.** 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

$$(\mathbf{I}_m - A)^2 = \mathbf{I}_m - 2A + A^2 = \mathbf{I}_m - A,$$

and so (a) holds. Let  $\lambda$  be an eigenvalue of  $A$  corresponding to the eigenvector  $\mathbf{x}$  so that  $A\mathbf{x} = \lambda\mathbf{x}$ . Then since  $A^2 = A$ , 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,  $A - \mathbf{I}_m$  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

$$\text{rank}(A) = r, \quad \text{rank}(A - \mathbf{I}_m) = m - r \quad (9.1)$$

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

$$\text{rank}(A) \geq r, \quad \text{rank}(A - \mathbf{I}_m) \geq m - r \quad (9.2)$$

But Corollary 2.12.1 gives

$$\text{rank}(A) + \text{rank}(\mathbf{I}_m - A) \leq \text{rank}\{A(\mathbf{I}_m - A)\} + m = \text{rank}\{(0)\} + m = m,$$

which together with (9.2) implies (9.1), 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. But 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  $\mathbf{I}_m$ .

If  $A$  is a diagonal matrix, that is,  $A = \text{diag}(a_1, \dots, a_m)$ , 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.

**Example 9.1.** 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 the example above 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 the following special case.

**Theorem 9.2.** 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  $\Lambda$  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  $\Lambda$ , that is, each eigenvalue of  $A$ , is 0 or 1.  $\square$

Our next result gives some conditions for the sum of two idempotent matrices and the product of two idempotent matrices to be idempotent.

**Theorem 9.3.** 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

$$AB = -BA \quad (9.3)$$

Premultiplication of (9.3) by  $B$  and postmultiplication by  $A$  yields the identity

$$(BA)^2 = -BA, \quad (9.4)$$

since  $A$  and  $B$  are idempotent. Similarly, premultiplying (9.3) by  $A$  and postmultiplying by  $B$ , we also find that

$$(AB)^2 = -AB \quad (9.5)$$

Thus, it follows from (9.4) and (9.5) that both  $-BA$  and  $-AB$  are idempotent matrices, and due to (9.3), so then is  $AB$ . Part (a) now follows since 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 9.2.** The conditions given for  $(A+B)$  to be idempotent are necessary and sufficient, while 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  $A$  and  $B$  are idempotent. In addition,  $AB = A$ , so that  $AB$  is also idempotent. However,  $AB \neq BA$  since  $BA = B$ .

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. The first result gives some restrictions on the elements of a symmetric idempotent matrix.

**Theorem 9.4.** 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

$$a_{ii} = (A)_{ii} = (A^2)_{ii} = (A'A)_{ii} = (A')_i \cdot (A)_{\cdot i} = \sum_{j=1}^m a_{ji}^2, \quad (9.6)$$

which clearly must be nonnegative. In addition, from (9.6) we have

$$a_{ii} = a_{ii}^2 + \sum_{j \neq i} a_{ji}^2,$$

so that  $a_{ii} \geq a_{ii}^2$  and, thus (b) must hold. If  $a_{ii} = 0$  or  $a_{ii} = 1$ , then  $a_{ii} = a_{ii}^2$  and so we must have

$$\sum_{j \neq i} a_{ji}^2 = 0,$$

which, along with the symmetry of  $A$ , establishes (c).  $\square$

The following theorem 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 9.5.** 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. But the identity  $A^{i+1} = A^i$  implies that  $\lambda_j^{i+1} = \lambda_j^i$ , for  $j = 1, \dots, m$ , so each  $\lambda_j$  must be either 0 or 1. The result now follows from Theorem 9.2.  $\square$

### 3. COCHRAN'S THEOREM

The following result, sometimes referred to as Cochran's Theorem [Cochran (1934)], will be very useful in establishing the independence of several different quadratic forms in the same normal variables.

**Theorem 9.6.** 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, there exists an orthogonal matrix  $P$  such that

$$P'A_hP = \text{diag}(I_r, (0))$$

For  $j \neq h$ , define  $B_j = P'A_jP$ , and note that

$$\begin{aligned} I_m &= P'I_mP = P'\left(\sum_{j=1}^k A_j\right)P = \left(\sum_{j=1}^k P'A_jP\right) \\ &= \text{diag}(I_r, (0)) + \sum_{j \neq h} B_j, \end{aligned}$$

or, equivalently,

$$\sum_{j \neq h} B_j = \text{diag}((0), I_{m-r})$$

In particular, for  $l = 1, \dots, r$ ,

$$\sum_{j \neq h} (B_j)_{ll} = 0$$

But clearly,  $B_j$  is symmetric and idempotent since  $A_j$  is, and so, from Theorem 9.4(a), its diagonal elements are nonnegative. Thus, we must have  $(B_j)_{ll} = 0$  for each  $l = 1, \dots, r$ , and this, along with Theorem 9.4(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  $j \neq h$ ,

$$P'A_hA_jP = (P'A_hP)(P'A_jP) = \text{diag}(I_r, (0)) \text{diag}((0), C_j) = (0),$$

which can be true only if  $A_hA_j = (0)$ , since  $P$  is nonsingular. Our proof is now complete, since  $h$  was arbitrary.  $\square$

Our next result is an extension of Cochran's Theorem.

**Theorem 9.7.** 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, there exists an orthogonal matrix  $P$  such that

$$P'AP = P'(A_1 + \dots + A_k)P = \text{diag}(I_r, (0)), \quad (9.7)$$

where  $r = \text{rank}(A)$ . Let  $B_i = P'A_iP$  for  $i = 1, \dots, k$ , and note that  $B_i$  is symmetric and idempotent. Thus, it follows from (9.7) and Theorem 9.4 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. But (9.7) also implies that

$$C_1 + \dots + C_k = I_r$$

Consequently,  $C_1, \dots, C_k$  satisfy the conditions of Theorem 9.6 and so  $C_i C_j = (0)$  for every  $i \neq j$ . From this 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, since

$$\begin{aligned} A^2 &= \left( \sum_{i=1}^k A_i \right)^2 = \sum_{i=1}^k \sum_{j=1}^k A_i A_j = \sum_{i=1}^k A_i^2 + \sum_{i \neq j} \sum A_i A_j \\ &= \sum_{i=1}^k A_i = A \end{aligned}$$

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 4.16, the matrices  $A_1, \dots, A_k$  can be simultaneously diagonalized; that is, there exists an orthogonal matrix  $Q$  such that

$$Q'A_iQ = D_i,$$

where each of the matrices  $D_1, \dots, D_k$  is diagonal. Further,

$$D_i D_j = Q'A_i Q Q'A_j Q = Q'A_i A_j Q = Q'(0)Q = (0), \quad (9.8)$$

for every  $i \neq j$ . Now since  $A$  is symmetric and idempotent so also is the diagonal matrix

$$Q'AQ = D_1 + \dots + D_k$$

As a result, each diagonal element of  $Q'AQ$  must be either 0 or 1, and due to (9.8) 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  $A_i = QD_iQ'$ . This completes the proof.  $\square$

Suppose that the three conditions given in Theorem 9.7 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 9.7 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 5.24 that  $F(F'HF)^-F'$  is a generalized inverse of  $H$  for any generalized inverse  $(F'HF)^-$ , of  $F'HF$ . But since  $A$  is idempotent,  $AI_mA = A$  and hence  $I_m$  is a generalized inverse of  $A = F'HF$ . Thus,  $FF'$  is a generalized inverse of  $H$ , yielding the equation

$$HFF'H = H,$$

which in partitioned form is

$$\begin{bmatrix} A_1^2 & A_1A_2 & \dots & A_1A_k \\ A_2A_1 & A_2^2 & \dots & A_2A_k \\ \vdots & \vdots & \ddots & \vdots \\ A_kA_1 & A_kA_2 & \dots & A_k^2 \end{bmatrix} = \begin{bmatrix} A_1 & (0) & \dots & (0) \\ (0) & A_2 & \dots & (0) \\ \vdots & \vdots & \ddots & \vdots \\ (0) & (0) & \dots & A_k \end{bmatrix}$$

This immediately gives conditions (a) and (c). The following result summarizes the relationship among these four conditions.

**Corollary 9.7.1.** 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)  $A_i A_j = (0)$ , 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.

#### 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 our first theorem to determine when the quadratic form  $x'Ax$  has a chi-squared distribution if the components of  $x$  are independently distributed, each having the  $N(0, 1)$  distribution.

**Theorem 9.8.** Let  $x \sim N_m(\mathbf{0}, I_m)$ , and suppose that the  $m \times m$  matrix  $A$  is symmetric, idempotent, and has rank  $r$ . Then  $x'Ax \sim \chi_r^2$ .

*Proof.* Since  $A$  is symmetric and idempotent, there exists an orthogonal matrix  $P$  such that

$$A = PDP',$$

where  $D = \text{diag}(I_r, (0))$ . Let  $z = P'x$  and note that since  $x \sim N_m(\mathbf{0}, I_m)$ ,

$$\begin{aligned} E(z) &= E(P'x) = P'E(x) = P'\mathbf{0} = \mathbf{0}, \\ \text{var}(z) &= \text{var}(P'x) = P'\{\text{var}(x)\}P = P'I_mP = P'P = I_m, \end{aligned}$$

and so  $z \sim N_m(\mathbf{0}, I_m)$ ; that is, the components of  $z$  are, like the components of  $x$ , independent standard normal random variables. Now due to 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 the result then follows. □

The result above is a special case of the next theorem in which the multivariate normal distribution has a general nonsingular covariance matrix.

**Theorem 9.9.** Let  $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  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  $\Omega$  is positive definite, there exists a nonsingular matrix  $T$  satisfying  $\Omega = TT'$ . If we define  $\mathbf{z} = T^{-1}\mathbf{x}$ , then  $E(\mathbf{z}) = T^{-1}E(\mathbf{x}) = \mathbf{0}$ , and

$$\text{var}(\mathbf{z}) = \text{var}(T^{-1}\mathbf{x}) = T^{-1}\text{var}(\mathbf{x})T'^{-1} = T^{-1}(TT')T'^{-1} = I_m,$$

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  $\mathbf{z}$  since

$$\mathbf{x}'A\mathbf{x} = \mathbf{x}'T'^{-1}T'ATT^{-1}\mathbf{x} = \mathbf{z}'T'AT\mathbf{z}$$

All that remains is to show that  $T'AT$  satisfies the conditions of the previous theorem. Clearly,  $T'AT$  is symmetric, since  $A$  is, and idempotent since

$$(T'AT)^2 = T'ATT'AT = T'A\Omega AT = T'AT,$$

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  $\Omega$  is nonsingular. Finally, since  $T'AT$  and  $A\Omega$  are idempotent, we have

$$\text{rank}(T'AT) = \text{tr}(T'AT) = \text{tr}(ATT') = \text{tr}(A\Omega) = \text{rank}(A\Omega) = r,$$

and so the proof is complete. □

It is not uncommon to have a quadratic form in a vector that has a singular multivariate normal distribution. Our next result generalizes the previous theorem to this situation.

**Theorem 9.10.** Let  $x \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  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  $x' A x \sim \chi_r^2$ .

*Proof.* Let  $n = \text{rank}(\Omega)$ , where  $n < m$ . Then there exists an  $m \times m$  orthogonal matrix  $P = [P_1 \ P_2]$  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  $\Lambda$  is an  $n \times n$  nonsingular diagonal matrix. Define

$$z = \begin{bmatrix} z_1 \\ z_2 \end{bmatrix} = \begin{bmatrix} P_1' x \\ P_2' x \end{bmatrix} = P' x,$$

and note that since  $P' \mathbf{0} = \mathbf{0}$  and  $P' \Omega P = \text{diag}(\Lambda, (0))$ ,  $z \sim N_m(\mathbf{0}, \text{diag}(\Lambda, (0)))$ . But this means that  $z = (z_1', \mathbf{0}')'$ , where  $z_1$  has the nonsingular distribution  $N_n(\mathbf{0}, \Lambda)$ . Now

$$x' A x = x' P P' A P P' x = z' P' A P z = z_1' P_1' A P_1 z_1,$$

and so the proof will be complete if we can show that the symmetric matrix  $P_1' A P_1$  satisfies the conditions of the previous theorem, namely, that  $P_1' A P_1 \Lambda$  is idempotent and  $\text{rank}(P_1' A P_1 \Lambda) = r$ . Since  $\Omega A \Omega A \Omega = \Omega A \Omega$ , we have

$$\begin{aligned} (\Lambda^{1/2} P_1' A P_1 \Lambda^{1/2})^3 &= \Lambda^{1/2} P_1' A (P_1 \Lambda P_1') A (P_1 \Lambda P_1') A P_1 \Lambda^{1/2} \\ &= \Lambda^{1/2} P_1' A \Omega A \Omega A P_1 \Lambda^{1/2} = \Lambda^{1/2} P_1' A \Omega A P_1 \Lambda^{1/2} \\ &= \Lambda^{1/2} P_1' A (P_1 \Lambda P_1') A P_1 \Lambda^{1/2} = (\Lambda^{1/2} P_1' A P_1 \Lambda^{1/2})^2, \end{aligned}$$

and so the idempotency of  $\Lambda^{1/2} P_1' A P_1 \Lambda^{1/2}$  follows from Theorem 9.5. However, this also establishes the idempotency of  $P_1' A P_1 \Lambda$  since  $\Lambda$  is nonsingular. Its rank is  $r$  since

$$\text{rank}(P_1' A P_1 \Lambda) = \text{tr}(P_1' A P_1 \Lambda) = \text{tr}(A P_1 \Lambda P_1') = \text{tr}(A \Omega) = r \quad \square$$

To this point, all of our results have dealt with normal distributions having the zero mean vector. In some applications, such as the determination of non-null distributions in hypothesis testing situations, we encounter quadratic forms in normal 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 9.9, is left to the reader. It utilizes 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 9.11.** Let  $x \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  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  $x'Ax \sim \chi_r^2(\lambda)$ , where  $\lambda = \frac{1}{2} \boldsymbol{\mu}'A\boldsymbol{\mu}$ .

**Theorem 9.12.** Let  $x \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  is positive semidefinite of rank  $n$ , and suppose that  $A$  is an  $m \times m$  symmetric matrix. Then  $x'Ax \sim \chi_r^2(\lambda)$ , 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  $\Lambda$  be defined as in the proof of Theorem 9.10 so that  $\Omega = P_1 \Lambda P_1'$ . Put  $C = [P_1 \Lambda^{-1/2} \quad P_2]$  and note that

$$\mathbf{z} = \begin{bmatrix} \mathbf{z}_1 \\ \mathbf{z}_2 \end{bmatrix} = \begin{bmatrix} \Lambda^{-1/2} P_1' \mathbf{x} \\ P_2' \mathbf{x} \end{bmatrix} = C' \mathbf{x} \sim N_m \left( \begin{bmatrix} \Lambda^{-1/2} P_1' \boldsymbol{\mu} \\ P_2' \boldsymbol{\mu} \end{bmatrix}, \begin{bmatrix} I_n & (0) \\ (0) & (0) \end{bmatrix} \right)$$

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 since  $C'^{-1} = [P_1 \Lambda^{1/2} \quad P_2]$ , we find that

$$\begin{aligned}
x'Ax &= x'CC^{-1}AC^{-1}C'x = z'C^{-1}AC^{-1}z \\
&= [z_1' \quad \mu'P_2] \begin{bmatrix} \Lambda^{1/2}P_1'AP_1\Lambda^{1/2} & \Lambda^{1/2}P_1'AP_2 \\ P_2'AP_1\Lambda^{1/2} & P_2'AP_2 \end{bmatrix} \begin{bmatrix} z_1 \\ P_2'\mu \end{bmatrix} \\
&= z_1'\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}z_1 + \mu'P_2P_2'AP_2P_2'\mu \\
&\quad + 2\mu'P_2P_2'AP_1\Lambda^{1/2}z_1 \tag{9.9}
\end{aligned}$$

But conditions (a)–(c) imply the identities

- (i)  $P_1'A\Omega AP_1 = P_1'AP_1$ ,
- (ii)  $\mu'P_2P_2'A\Omega AP_1 = \mu'P_2P_2'AP_1$ ,
- (iii)  $\mu'P_2P_2'A\Omega A\Omega AP_2P_2'\mu = \mu'P_2P_2'A\Omega AP_2P_2'\mu = \mu'P_2P_2'AP_2P_2'\mu$ ;

in particular, (a) implies (i), (b) and (i) imply (ii), while (iii) follows from (c), (i) and (ii). Utilizing these identities in (9.9), we obtain

$$\begin{aligned}
x'Ax &= z_1'\Lambda^{1/2}P_1'AP_1\Lambda^{1/2}z_1 + \mu'P_2P_2'A\Omega A\Omega AP_2P_2'\mu \\
&\quad + 2\mu'P_2P_2'A\Omega AP_1\Lambda^{1/2}z_1 \\
&= (z_1 + \Lambda^{1/2}P_1'AP_2P_2'\mu)' \Lambda^{1/2}P_1'AP_1\Lambda^{1/2} (z_1 + \Lambda^{1/2}P_1'AP_2P_2'\mu) \\
&= w'A_*w.
\end{aligned}$$

Now,  $w = (z_1 + \Lambda^{1/2}P_1'AP_2P_2'\mu) \sim N_n(\theta, I_n)$ , where

$$\theta = \Lambda^{-1/2}P_1'\mu + \Lambda^{1/2}P_1'AP_2P_2'\mu,$$

and, since  $A_* = \Lambda^{1/2}P_1'AP_1\Lambda^{1/2}$  is idempotent, a consequence of (i), we may apply Theorem 9.11; that is,  $w'A_*w \sim \chi_r^2(\lambda)$ , 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

$$\begin{aligned}
\lambda &= \frac{1}{2} \theta'A_*\theta = \frac{1}{2} (\Lambda^{-1/2}P_1'\mu + \Lambda^{1/2}P_1'AP_2P_2'\mu)' \\
&\quad \times \Lambda^{1/2}P_1'AP_1\Lambda^{1/2}(\Lambda^{-1/2}P_1'\mu + \Lambda^{1/2}P_1'AP_2P_2'\mu) \\
&= \frac{1}{2} (\mu'P_1P_1'AP_1P_1'\mu + \mu'P_2P_2'A\Omega A\Omega AP_2P_2'\mu + 2\mu'P_1P_1'A\Omega AP_2P_2'\mu)
\end{aligned}$$

$$\begin{aligned}
 &= \frac{1}{2} (\boldsymbol{\mu}' P_1 P_1' A P_1 P_1' \boldsymbol{\mu} + \boldsymbol{\mu}' P_2 P_2' A P_2 P_2' \boldsymbol{\mu} + 2 \boldsymbol{\mu}' P_1 P_1' A P_2 P_2' \boldsymbol{\mu}) \\
 &= \frac{1}{2} \boldsymbol{\mu}' (P_1 P_1' + P_2 P_2') A (P_1 P_1' + P_2 P_2') \boldsymbol{\mu} = \frac{1}{2} \boldsymbol{\mu}' A \boldsymbol{\mu}
 \end{aligned}$$

This completes the proof. □

A matrix  $A$  satisfying conditions (a), (b), and (c) of Theorem 9.12 is  $\Omega^+$ , the Moore–Penrose inverse of  $\Omega$ . That is, if  $x \sim N_m(\boldsymbol{\mu}, \Omega)$ , then  $x' \Omega^+ x$  will have a chi-squared distribution since the identity  $\Omega^+ \Omega \Omega^+ = \Omega^+$  ensures that conditions (a), (b), and (c) hold. The degrees of freedom  $r = \text{rank}(\Omega)$  since  $\text{rank}(\Omega^+ \Omega) = \text{rank}(\Omega)$ .

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. This is most easily proven using moment generating functions. For details on this, the interested reader is referred to Mathai and Provost (1992) or Searle (1971).

**Example 9.3.** Let  $x_1, \dots, x_n$  be a random sample from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ ; 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  $x = (x_1, \dots, x_n)'$  so that  $x \sim N_n(\mu \mathbf{1}_n, \sigma^2 \mathbf{I}_n)$ . Note that if the  $n \times n$  matrix  $A = (\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}_n') / \sigma^2$ , then

$$\begin{aligned}
 x' A x &= \frac{\{x' x - n^{-1} (\mathbf{1}_n' x)^2\}}{\sigma^2} = \left\{ \sum_{i=1}^n x_i^2 - n^{-1} \left( \sum_{i=1}^n x_i \right)^2 \right\} / \sigma^2 \\
 &= \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{\sigma^2} = \frac{(n-1)s^2}{\sigma^2} = t,
 \end{aligned}$$

and so  $t$  is a quadratic form in the random vector  $x$ . The matrix  $A(\sigma^2 \mathbf{I}_n) = \sigma^2 A$

is idempotent since

$$\begin{aligned}(\sigma^2 A)^2 &= (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n)^2 = I_n - 2n^{-1} \mathbf{1}_n \mathbf{1}'_n + n^{-2} \mathbf{1}_n \mathbf{1}'_n \mathbf{1}_n \mathbf{1}'_n \\ &= I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n = \sigma^2 A,\end{aligned}$$

and so by Theorem 9.11,  $t$  has a chi-squared distribution. This chi-squared distribution has  $n - 1$  degrees of freedom since

$$\text{tr}(\sigma^2 A) = \text{tr}(I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) = \text{tr}(I_n) - n^{-1} \text{tr}(\mathbf{1}_n \mathbf{1}'_n) = n - n^{-1} \mathbf{1}'_n \mathbf{1}_n = n - 1,$$

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  $t \sim \chi_{n-1}^2$ .

## 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 or not 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 9.13.** Let  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  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  $\Omega$  is positive definite, there exists a nonsingular matrix  $T$  such that  $\Omega = TT'$ . Define  $G = T'AT$  and  $H = T'BT$ , and note that if  $A\Omega B = (0)$ , then

$$GH = (T'AT)(T'BT) = T'A\Omega BT = T'(0)T = (0) \quad (9.10)$$

Consequently, due to 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 4.15, we know that there exists an orthogonal matrix  $P$  that simultaneously diagonalizes  $G$  and  $H$ ; that is, for some diagonal matrices  $C$  and  $D$ ,

$$P'GP = P'T'ATP = C, \quad P'HP = P'T'BTTP = D \quad (9.11)$$

But using (9.10) and (9.11), we find that

$$(0) = GH = PCP'PDP' = PCDP',$$

which only can be true 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  $y = P'T^{-1}x$ , then our two quadratic forms simplify as

$$x'Ax = x'T'^{-1}PP'T'ATPP'T^{-1}x = y'Cy = \sum_{i=1}^{m_1} c_i y_i^2,$$

and

$$x'Bx = x'T'^{-1}PP'T'BTTP'T^{-1}x = y'Dy = \sum_{i=m_1+1}^m d_i y_i^2;$$

that is, the first quadratic form is a function only of  $y_1, \dots, y_{m_1}$ , while the second quadratic form is a function of  $y_{m_1+1}, \dots, y_m$ . The result now follows from the independence of  $y_1, \dots, y_m$ , a consequence of the fact that  $y$  is normal and

$$\text{var}(y) = \text{var}(P'T^{-1}x) = P'T^{-1}\Omega T'^{-1}P = I_m \quad \square$$

**Example 9.4.** Suppose that  $x_1, \dots, x_k$  are independently distributed with  $x_i = (x_{i1}, \dots, x_{in})' \sim N_n(\mu \mathbf{1}_n, \sigma^2 I_n)$  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

$$\bar{x}_i = \sum_{j=1}^n \frac{x_{ij}}{n}, \quad \bar{\mathbf{x}} = \sum_{i=1}^k \frac{\bar{x}_i}{k}$$

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 7.4). Now  $t_1$  can be expressed as

$$t_1 = n \left\{ \sum_{i=1}^k \bar{x}_i^2 - k^{-1} \left( \sum_{i=1}^k \bar{x}_i \right)^2 \right\} = n \bar{\mathbf{x}}' (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \bar{\mathbf{x}},$$

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  $\Omega = \mathbf{I}_k \otimes \sigma^2 \mathbf{I}_n = \sigma^2 \mathbf{I}_{kn}$ , and  $\bar{\mathbf{x}} = n^{-1} (\mathbf{I}_k \otimes \mathbf{1}'_n) \mathbf{x}$ , so

$$\begin{aligned} t_1 &= n^{-1} \mathbf{x}' (\mathbf{I}_k \otimes \mathbf{1}_n) (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') (\mathbf{I}_k \otimes \mathbf{1}'_n) \mathbf{x} \\ &= n^{-1} \mathbf{x}' \{ (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}'_n \} \mathbf{x} = \mathbf{x}' A_1 \mathbf{x}, \end{aligned}$$

where  $A_1 = n^{-1} \{ (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}'_n \}$ . Since  $(\mathbf{1}_n \mathbf{1}'_n)^2 = n \mathbf{1}_n \mathbf{1}'_n$  and  $(\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k')^2 = (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k')$ , we find that  $A_1$  is idempotent and hence so is  $(A_1/\sigma^2)\Omega$ . Thus, by Theorem 9.11,  $\mathbf{x}' (A_1/\sigma^2) \mathbf{x} = t_1/\sigma^2$  has a chi-squared distribution. This distribution is central since  $\lambda = \frac{1}{2} \boldsymbol{\mu}' A_1 \boldsymbol{\mu} / \sigma^2 = 0$ , which follows from the fact that

$$\begin{aligned} \{ (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}'_n \} (\mathbf{1}_k \otimes \mu \mathbf{1}_n) &= n \mu \{ (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \} \\ &= n \mu \{ (\mathbf{1}_k - \mathbf{1}_k) \otimes \mathbf{1}_n \} = \mathbf{0}, \end{aligned}$$

while its degrees of freedom is given by

$$\begin{aligned} r_1 &= \text{tr} \{ (A_1/\sigma^2) \Omega \} = \text{tr}(A_1) = n^{-1} \text{tr} \{ (\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \otimes \mathbf{1}_n \mathbf{1}'_n \} \\ &= n^{-1} \text{tr}(\mathbf{I}_k - k^{-1} \mathbf{1}_k \mathbf{1}_k') \text{tr}(\mathbf{1}_n \mathbf{1}'_n) = n^{-1} (k-1)n = k-1 \end{aligned}$$

Turning to  $t_2$ , observe that it can be written as

$$\begin{aligned} t_2 &= \sum_{i=1}^k \left\{ \sum_{j=1}^n x_{ij}^2 - n^{-1} \left( \sum_{j=1}^n x_{ij} \right)^2 \right\} = \sum_{i=1}^k \mathbf{x}'_i (\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) \mathbf{x}_i \\ &= \mathbf{x}' \{ \mathbf{I}_k \otimes (\mathbf{I}_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) \} \mathbf{x} = \mathbf{x}' A_2 \mathbf{x}, \end{aligned}$$

where  $A_2 = I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n)$ . Clearly,  $A_2$  is idempotent since  $(I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n)$  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,  $t_2/\sigma^2 \sim \chi^2_{k(n-1)}$  since

$$\begin{aligned} \text{tr}\{(A_2/\sigma^2)\Omega\} &= \text{tr}(A_2) = \text{tr}\{I_k \otimes (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n)\} \\ &= \text{tr}(I_k)\text{tr}(I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) = k(n - 1), \end{aligned}$$

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 9.13. This simply involves verifying that  $(A_1/\sigma^2)\Omega(A_2/\sigma^2) = A_1A_2/\sigma^2 = (0)$ , which is an immediate consequence of the fact that

$$\mathbf{1}_n \mathbf{1}'_n (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) = (0)$$

**Example 9.5.** 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 since 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

$$t_1 = (\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1)'(\mathbf{y} - X_1\hat{\boldsymbol{\beta}}_1) = \mathbf{y}'(I_N - X_1(X_1'X_1)^{-1}X_1')\mathbf{y},$$

to the sum of squared errors for the complete model, which is given by

$$t_2 = (\mathbf{y} - X\hat{\boldsymbol{\beta}})'(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = \mathbf{y}'(I_N - X(X'X)^{-1}X')\mathbf{y}$$

Now if  $\boldsymbol{\epsilon} \sim N_N(\mathbf{0}, \sigma^2 I)$ , then  $\mathbf{y} \sim N_N(X\boldsymbol{\beta}, \sigma^2 I)$ . Thus, by applying Theorem 9.11 and using the fact that  $X(X'X)^{-1}X'X_1 = X_1$ , we find that  $(t_1 - t_2)/\sigma^2$  is chi-squared since

$$\left\{ \frac{(X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1')}{\sigma^2} \right\} \{\sigma^2 I\} \left\{ \frac{(X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1')}{\sigma^2} \right\} \\ = \left\{ \frac{(X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1')}{\sigma^2} \right\}$$

In particular, if  $\beta_2 = \mathbf{0}$ ,  $(t_1 - t_2)/\sigma^2 \sim \chi_{m_2}^2$ , since

$$\text{tr}\{X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1'\} = \text{tr}\{X(X'X)^{-1}X'\} - \text{tr}\{X_1(X_1'X_1)^{-1}X_1'\} \\ = r - (r - m_2) = m_2,$$

and

$$\beta_1' X_1' \left\{ \frac{(X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1')}{\sigma^2} \right\} X_1 \beta_1 = \frac{(\beta_1' X_1' X_1 \beta_1 - \beta_1' X_1' X_1 \beta_1)}{\sigma^2} = 0$$

By a similar application of Theorem 9.11, we observe that  $t_2/\sigma^2 \sim \chi_{N-r}^2$ . In addition, it follows from Theorem 9.13 that  $(t_1 - t_2)/\sigma^2$  and  $t_2/\sigma^2$  are independently distributed since

$$\left\{ \frac{(X(X'X)^{-1}X' - X_1(X_1'X_1)^{-1}X_1')}{\sigma^2} \right\} \{\sigma^2 I\} \left\{ \frac{(I_N - X(X'X)^{-1}X')}{\sigma^2} \right\} = 0$$

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 of the next result, which is very similar to the proof of Theorem 9.13, is left to the reader as an exercise.

**Theorem 9.14.** Let  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  is positive definite, and suppose that  $A$  is an  $m \times m$  symmetric matrix while  $B$  is an  $n \times m$  matrix. If  $B\Omega A = (\mathbf{0})$ , then  $\mathbf{x}'A\mathbf{x}$  and  $B\mathbf{x}$  are independently distributed.

**Example 9.6.** Suppose that we have a random sample  $x_1, \dots, x_n$  from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ . In Example 9.3, 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 9.14 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 9.3, we saw that  $s^2$  is a scalar multiple of the quadratic form

$$\mathbf{x}'(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n')\mathbf{x},$$

where  $\mathbf{x} = (x_1, \dots, x_n)' \sim N_n(\mu\mathbf{1}_n, \sigma^2\mathbf{I}_n)$ . 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

$$\mathbf{1}_n'(\sigma^2\mathbf{I}_n)(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n') = \sigma^2(\mathbf{1}_n' - n^{-1}\mathbf{1}_n'\mathbf{1}_n\mathbf{1}_n') = \sigma^2(\mathbf{1}_n' - \mathbf{1}_n') = \mathbf{0}'$$

When  $\Omega$  is positive semidefinite, the condition  $A\Omega B = (0)$ , given in Theorem 9.13, will still guarantee that the two quadratic forms  $\mathbf{x}'A\mathbf{x}$  and  $\mathbf{x}'B\mathbf{x}$  are independently distributed. Likewise, when  $\Omega$  is positive semidefinite, the condition  $B\Omega A = (0)$ , given in Theorem 9.14, 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 the following two theorems. The proofs are left as exercises.

**Theorem 9.15.** Let  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  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 9.16.** Let  $x \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  is positive semidefinite, and suppose that  $A$  is an  $m \times m$  symmetric matrix while  $B$  is an  $n \times m$  matrix. If  $B\Omega A\Omega = (0)$  and  $B\Omega A\boldsymbol{\mu} = \mathbf{0}$ , then  $x'Ax$  and  $Bx$  are independently distributed.

Our final result can be helpful in establishing that several quadratic forms in the same normal random vector are independently distributed, each having a chi-squared distribution.

**Theorem 9.17.** Let  $x \sim N_m(\boldsymbol{\mu}, \Omega)$ , where  $\Omega$  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)  $x'A_i x \sim \chi_{r_i}^2(\frac{1}{2}\boldsymbol{\mu}'A_i\boldsymbol{\mu})$ ,
- (ii)  $x'Ax \sim \chi_r^2(\frac{1}{2}\boldsymbol{\mu}'A\boldsymbol{\mu})$ ,
- (iii)  $x'A_1x, \dots, x'A_kx$  are independently distributed.

*Proof.* Since  $\Omega$  is positive definite, there exists a nonsingular matrix  $T$  satisfying  $\Omega = TT'$ , 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 9.7.1, 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) above hold. Now using Theorem 9.11, (a) implies (i) and (b) implies (ii), while Theorem 9.13, along with (c), guarantees that (iii) holds.  $\square$

## 6. EXPECTED VALUES OF QUADRATIC FORMS

When a quadratic form satisfies the conditions given in the theorems of Section 6.4, 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  $x$  has an arbitrary distribution. The expressions we obtain involve the matrix of second moments of  $x$ ,  $E(xx')$ , and the matrix of fourth moments  $E(xx' \otimes xx')$ .

**Theorem 9.18.** Let  $x$  be an  $m \times 1$  random vector having finite fourth moments, so that both  $E(xx')$  and  $E(xx' \otimes xx')$  exist. Denote the mean vector and covariance matrix of  $x$  by  $\mu$  and  $\Omega$ . If  $A$  and  $B$  are  $m \times m$  symmetric matrices, then

- (a)  $E(x'Ax) = \text{tr}\{AE(xx')\} = \text{tr}(A\Omega) + \mu'A\mu,$
- (b)  $\text{var}(x'Ax) = \text{tr}\{(A \otimes A)E(xx' \otimes xx')\} - [\text{tr}(A\Omega) + \mu'A\mu]^2,$
- (c)  $\text{cov}(x'Ax, x'Bx) = \text{tr}\{(A \otimes B)E(xx' \otimes xx')\} - [\text{tr}(A\Omega) + \mu'A\mu][\text{tr}(B\Omega) + \mu'B\mu].$

*Proof.* The covariance matrix  $\Omega$  is defined by

$$\Omega = E\{(x - \mu)(x - \mu)'\} = E(xx') - \mu\mu'$$

so that  $E(xx') = \Omega + \mu\mu'$ . Since  $x'Ax$  is a scalar, we have

$$\begin{aligned} E(x'Ax) &= E\{\text{tr}(x'Ax)\} = E\{\text{tr}(Axx')\} = \text{tr}\{AE(xx')\} = \text{tr}\{A(\Omega + \mu\mu')\} \\ &= \text{tr}(A\Omega) + \text{tr}(A\mu\mu') = \text{tr}(A\Omega) + \mu'A\mu, \end{aligned}$$

and so (a) holds. Part (b) will follow from (c) by taking  $B = A$ . To prove (c), note that

$$\begin{aligned} E(x'Ax x'Bx) &= E[\text{tr}\{(x' \otimes x')(A \otimes B)(x \otimes x)\}] \\ &= E[\text{tr}\{(A \otimes B)(x \otimes x)(x' \otimes x')\}] \\ &= \text{tr}\{(A \otimes B)E(xx' \otimes xx')\} \end{aligned}$$

Then use this, along with part (a), in the equation

$$\text{cov}(x'Ax, x'Bx) = E(x'Ax x'Bx) - E(x'Ax)E(x'Bx) \quad \square$$

When  $x$  has a normal distribution, the expressions for variances and covariances, as well as 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 7, plays a crucial role in obtaining some of these matrix expressions. We will also make use of the  $m \times m$  matrix  $T_{ij}$  defined by

$$T_{ij} = E_{ij} + E_{ji} = e_i e_j' + e_j e_i';$$

that is, all of the 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 the following result.

**Theorem 9.19.** If  $z \sim N_m(\mathbf{0}, I_m)$  and  $c$  is a vector of constants, then

- (a)  $E(z \otimes z) = \text{vec}(I_m)$ ,
- (b)  $E(cz' \otimes zz') = (0)$ ,  $E(zc' \otimes zz') = (0)$ ,  $E(zz' \otimes cz') = (0)$ ,  $E(zz' \otimes zc') = (0)$ ,
- (c)  $E(zz' \otimes zz') = 2N_m + \text{vec}(I_m)\{\text{vec}(I_m)\}'$ ,
- (d)  $\text{var}(z \otimes z) = 2N_m$ .

*Proof.* Since  $E(z) = \mathbf{0}$ ,  $I_m = \text{var}(z) = E(zz')$  and so

$$E(z \otimes z) = E\{\text{vec}(zz')\} = \text{vec}\{E(zz')\} = \text{vec}(I_m)$$

It is easily verified using the standard normal moment generating function that

$$E(z_i^3) = 0, \quad E(z_i^4) = 3 \quad (9.12)$$

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  $z$  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,

$$E(z_j z_k z_l) = E(z_j^2)E(z_l) = 1 \cdot 0 = 0$$

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 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  $\delta_{ij}$  is the  $(i, j)$ th element of  $I_m$ . Thus,

$$\begin{aligned} E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}') &= E\left\{\left(\sum_{i=1}^m \sum_{j=1}^m E_{ij}z_i z_j\right) \otimes \mathbf{z}\mathbf{z}'\right\} = \sum_{i=1}^m \sum_{j=1}^m \{E_{ij} \otimes E(z_i z_j \mathbf{z}\mathbf{z}')\} \\ &= \sum_{i=1}^m \sum_{j=1}^m \{E_{ij} \otimes (T_{ij} + \delta_{ij}I_m)\} \\ &= \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes T_{ij}) + \sum_{i=1}^m \sum_{j=1}^m (\delta_{ij}E_{ij} \otimes I_m) \end{aligned}$$

The third result now follows since

$$\begin{aligned} \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes T_{ij}) &= \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes E_{ji}) + \sum_{i=1}^m \sum_{j=1}^m (E_{ij} \otimes E_{ij}) \\ &= K_{mm} + \left\{ \sum_{i=1}^m (e_i \otimes e_i) \right\} \left\{ \sum_{j=1}^m (e'_j \otimes e'_j) \right\} \\ &= K_{mm} + \left\{ \sum_{i=1}^m \text{vec}(e_i e'_i) \right\} \left\{ \sum_{j=1}^m \{\text{vec}(e_j e'_j)\}' \right\} \\ &= K_{mm} + \text{vec}(I_m) \{\text{vec}(I_m)\}', \\ \sum_{i=1}^m \sum_{j=1}^m (\delta_{ij}E_{ij} \otimes I_m) &= \left( \sum_{i=1}^m E_{ii} \right) \otimes I_m = I_m \otimes I_m = I_{m^2}, \end{aligned}$$

and  $I_{m^2} + K_{mm} = 2N_m$ . Finally, (d) is an immediate consequence of (a) and (c). □

The next result generalizes the results of Theorem 9.19 to a multivariate normal distribution having a general positive definite covariance matrix.

**Theorem 9.20.** Let  $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  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)  $E(\mathbf{c}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') = (0)$ ,  $E(\mathbf{x}\mathbf{c}' \otimes \mathbf{x}\mathbf{x}') = (0)$ ,  $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{c}\mathbf{x}') = (0)$ ,  $E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{c}') = (0)$ ,

$$(c) E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') = 2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)\}',$$

$$(d) \text{var}(\mathbf{x} \otimes \mathbf{x}) = 2N_m(\Omega \otimes \Omega).$$

*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 above are consequences of Theorem 9.19 since

$$\begin{aligned} E(\mathbf{x} \otimes \mathbf{x}) &= (T \otimes T)E(\mathbf{z} \otimes \mathbf{z}) = (T \otimes T)\text{vec}(I_m) = \text{vec}(TT') = \text{vec}(\Omega), \\ E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') &= (I_m \otimes T)E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}')(T' \otimes T') = (0), \end{aligned}$$

and

$$\begin{aligned} E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') &= (T \otimes T)E(\mathbf{z}\mathbf{z}' \otimes \mathbf{z}\mathbf{z}')(T' \otimes T') \\ &= (T \otimes T)(2N_m + \text{vec}(I_m)\{\text{vec}(I_m)\}')(T' \otimes T') \\ &= 2(T \otimes T)N_m(T' \otimes T') + (T \otimes T)\text{vec}(I_m)\{\text{vec}(I_m)\}'(T' \otimes T') \\ &= 2N_m(T \otimes T)(T' \otimes T') + \text{vec}(TT')\{\text{vec}(TT')\}' \\ &= 2N_m(\Omega \otimes \Omega) + \text{vec}(\Omega)\{\text{vec}(\Omega)\}' \quad \square \end{aligned}$$

We are now ready to obtain simplified expressions for the variance and covariance of quadratic forms in normal variates.

**Theorem 9.21.** Let  $A$  and  $B$  be  $m \times m$  symmetric matrices and suppose that  $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  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 9.20, 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 7.8, while

$$\text{tr}\{(A \otimes B)K_{mm}(\Omega \otimes \Omega)\} = \text{tr}\{(A\Omega \otimes B\Omega)K_{mm}\} = \text{tr}(A\Omega B\Omega)$$

follows from Theorem 7.31. Using the symmetry of  $A$  and  $\Omega$  along with Theorems 7.15 and 7.16, the last term in  $E\{x'Ax x'Bx\}$  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). To prove (b) we use the definition of covariance and Theorem 9.18(a) to get

$$\text{cov}(x'Ax, x'Bx) = E\{x'Ax x'Bx\} - E(x'Ax)E(x'Bx) = 2 \text{tr}(A\Omega B\Omega) \quad \square$$

The formulas given in the previous theorem become somewhat more complicated when the normal distribution has a nonnull mean vector. These formulas are given in the following theorem.

**Theorem 9.22.** Let  $A$  and  $B$  be symmetric  $m \times m$  matrices and suppose that  $x \sim N_m(\mu, \Omega)$ , where  $\Omega$  is positive definite. Then

- (a)  $E\{x'Ax x'Bx\} = \text{tr}(A\Omega)\text{tr}(B\Omega) + 2 \text{tr}(A\Omega B\Omega) + \text{tr}(A\Omega)\mu'B\mu + 4\mu'A\Omega B\mu + \mu'A\mu \text{tr}(B\Omega) + \mu'A\mu\mu'B\mu,$
- (b)  $\text{cov}(x'Ax, x'Bx) = 2 \text{tr}(A\Omega B\Omega) + 4\mu'A\Omega B\mu,$
- (c)  $\text{var}(x'Ax) = 2 \text{tr}\{(A\Omega)^2\} + 4\mu'A\Omega A\mu.$

*Proof.* Again (c) is a special case of (b), so we only need to prove (a) and (b). We can write  $x = y + \mu$ , where  $y \sim N_m(\mathbf{0}, \Omega)$  and, consequently,

$$\begin{aligned} E\{x'Ax x'Bx\} &= E\{(y + \mu)'A(y + \mu)(y + \mu)'B(y + \mu)\} \\ &= E\{(y'Ay + 2\mu'Ay + \mu'A\mu)(y'By + 2\mu'By + \mu'B\mu)\} \\ &= E\{y'Ay y'By\} + 2E\{y'Ay \mu'By\} + E(y'Ay)\mu'B\mu \\ &\quad + 2E\{\mu'Ay y'By\} + 4E(\mu'Ay \mu'By) + 2E(\mu'Ay)\mu'B\mu \\ &\quad + \mu'A\mu E(y'By) + 2\mu'A\mu E(\mu'By) + \mu'A\mu \mu'B\mu \end{aligned}$$

The sixth and eighth terms in this last expression are zero since  $E(y) = \mathbf{0}$ , while it follows from Theorem 9.20(b) that the second and fourth terms are zero. To simplify the fifth term note that

$$\begin{aligned} E(\boldsymbol{\mu}'Ay \boldsymbol{\mu}'By) &= E\{(\boldsymbol{\mu}'A \otimes \boldsymbol{\mu}'B)(y \otimes y)\} = (A\boldsymbol{\mu} \otimes B\boldsymbol{\mu})'E\{(y \otimes y)\} \\ &= \{\text{vec}(B\boldsymbol{\mu}\boldsymbol{\mu}'A)\}' \text{vec}(\Omega) = \text{tr}\{(B\boldsymbol{\mu}\boldsymbol{\mu}'A)'\Omega\} \\ &= \text{tr}(A\boldsymbol{\mu}\boldsymbol{\mu}'B\Omega) = \boldsymbol{\mu}'A\Omega B\boldsymbol{\mu} \end{aligned}$$

Thus, using this and Theorems 9.18(a) and 9.21(a), we find that

$$\begin{aligned} E\{\mathbf{x}'A\mathbf{x} \mathbf{x}'B\mathbf{x}\} &= \text{tr}(A\Omega) \text{tr}(B\Omega) + 2 \text{tr}(A\Omega B\Omega) + \text{tr}(A\Omega)\boldsymbol{\mu}'B\boldsymbol{\mu} + 4\boldsymbol{\mu}'A\Omega B\boldsymbol{\mu} \\ &\quad + \boldsymbol{\mu}'A\boldsymbol{\mu} \text{tr}(B\Omega) + \boldsymbol{\mu}'A\boldsymbol{\mu} \boldsymbol{\mu}'B\boldsymbol{\mu}, \end{aligned}$$

thereby proving (a); (b) then follows immediately from the definition of covariance and Theorem 9.18(a).  $\square$

**Example 9.7.** Let us return to the subject of Example 9.4, where we defined

$$A_1 = n^{-1} \{(\mathbf{I}_k - k^{-1}\mathbf{1}_k\mathbf{1}_k') \otimes \mathbf{1}_n\mathbf{1}_n'\}$$

and

$$A_2 = \mathbf{I}_k \otimes (\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n')$$

It was shown that if  $\mathbf{x} = (\mathbf{x}'_1, \dots, \mathbf{x}'_k)'\sim N_{kn}(\boldsymbol{\mu}, \Omega)$  with  $\boldsymbol{\mu} = \mathbf{1}_k \otimes \boldsymbol{\mu}\mathbf{1}_n$  and  $\Omega = \mathbf{I}_k \otimes \sigma^2\mathbf{I}_n$ , then  $t_1/\sigma^2 = \mathbf{x}'(A_1/\sigma^2)\mathbf{x} \sim \chi_{k-1}^2$ ,  $t_2/\sigma^2 = \mathbf{x}'(A_2/\sigma^2)\mathbf{x} \sim \chi_{k(n-1)}^2$ , independently. Since the mean of a chi-squared random variable equals its degrees of freedom, while the variance is two times the degrees of freedom, we can easily calculate the mean and variance of  $t_1$  and  $t_2$  without employing 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  $\mathbf{x}_i \sim N_n(\boldsymbol{\mu}\mathbf{1}_n, \sigma_i^2\mathbf{I}_n)$  so that  $\Omega = \text{var}(\mathbf{x}) = D \otimes \mathbf{I}_n$ , where  $D = \text{diag}(\sigma_1^2, \dots, \sigma_k^2)$ . It can be easily verified that, in this case,  $t_1/\sigma^2$  and  $t_2/\sigma^2$  no longer satisfy the conditions of Theorem 9.11 for chi-squaredness, but are still independently distributed. The mean and variance of  $t_1$  and  $t_2$  can be computed by using Theorems 9.18 and 9.22. For instance, the mean of  $t_2$  is given by

$$\begin{aligned} E(t_2) &= E(\mathbf{x}'A_2\mathbf{x}) = \text{tr}(A_2\Omega) + \boldsymbol{\mu}'A_2\boldsymbol{\mu} \\ &= \text{tr}\{(\mathbf{I}_k \otimes (\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n'))(D \otimes \mathbf{I}_n)\} \\ &\quad + \boldsymbol{\mu}^2(\mathbf{1}'_k \otimes \mathbf{1}'_n)\{\mathbf{I}_k \otimes (\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}_n')\}(\mathbf{1}_k \otimes \mathbf{1}_n) \end{aligned}$$



$$\begin{aligned}
 &= \text{tr}(D)\text{tr}(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n) + \mu^2(\mathbf{1}'_k\mathbf{1}_k)\{\mathbf{1}'_n(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n)\mathbf{1}_n\} \\
 &= (n - 1) \sum_{i=1}^k \sigma_i^2,
 \end{aligned}$$

while its variance is

$$\begin{aligned}
 \text{var}(t_2) &= \text{var}(\mathbf{x}'A_2\mathbf{x}) = 2 \text{tr}\{(A_2\Omega)^2\} + 4\mu'A_2\Omega A_2\mu \\
 &= 2 \text{tr}\{D^2 \otimes (\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n)\} \\
 &\quad + 4\mu^2(\mathbf{1}'_k \otimes \mathbf{1}'_n)\{D \otimes (\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n)\}(\mathbf{1}_k \otimes \mathbf{1}_n) \\
 &= 2 \text{tr}(D^2)\text{tr}\{(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n)\} + 4\mu^2(\mathbf{1}'_k D \mathbf{1}_k)\{\mathbf{1}'_n(\mathbf{I}_n - n^{-1}\mathbf{1}_n\mathbf{1}'_n)\mathbf{1}_n\} \\
 &= 2(n - 1) \sum_{i=1}^k \sigma_i^4
 \end{aligned}$$

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 using this 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, utilizing 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 below.

**Theorem 9.23.** 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

$$(a) \ E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}) = \text{tr}(A)\text{tr}(B)\text{tr}(C) + 2\{\text{tr}(A)\text{tr}(BC) + \text{tr}(B)\text{tr}(AC) + \text{tr}(C)\text{tr}(AB)\} + 8\text{tr}(ABC),$$

$$(b) \ E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x}\mathbf{x}'D\mathbf{x}) \\ = \text{tr}(A)\text{tr}(B)\text{tr}(C)\text{tr}(D) + 8\{\text{tr}(A)\text{tr}(BCD) + \text{tr}(B)\text{tr}(ACD) + \text{tr}(C)\text{tr}(ABD) + \text{tr}(D)\text{tr}(ABC)\} + 4\{\text{tr}(AB)\text{tr}(CD) + \text{tr}(AC)\text{tr}(BD) + \text{tr}(AD)\text{tr}(BC)\} + 2\{\text{tr}(A)\text{tr}(B)\text{tr}(CD) + \text{tr}(A)\text{tr}(C)\text{tr}(BD) + \text{tr}(A)\text{tr}(D)\text{tr}(BC) + \text{tr}(B)\text{tr}(C)\text{tr}(AD) + \text{tr}(B)\text{tr}(D)\text{tr}(AC) + \text{tr}(C)\text{tr}(D)\text{tr}(AB)\} + 16\{\text{tr}(ABCD) + \text{tr}(ABDC) + \text{tr}(ACBD)\}$$

If  $\mathbf{x} \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  is positive definite, then  $A$ ,  $B$ ,  $C$ , and  $D$  appearing in the right-hand side of the equations in Theorem 9.23 are replaced by  $A\Omega$ ,  $B\Omega$ ,  $C\Omega$ , and  $D\Omega$ .

An alternative approach to the calculation of moments of quadratic forms utilizes 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).

## 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 which 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  $\sigma_{jj}$  is the  $j$ th diagonal element of  $\Omega$ , so that the  $j$ th diagonal element of  $X'X$  has the distribution  $\sigma_{jj}\chi_n^2$ . The joint distribution of all of the elements of the  $m \times m$  matrix  $X'X$  is called the Wishart distribution with scale matrix  $\Omega$  and degrees of freedom  $n$ , and will be denoted by  $W_m(\Omega, n)$ . This Wishart distribution, like the chi-squared distribution  $\chi_n^2$ , is said to be central. More generally, if  $x_1, \dots, x_n$  are independent and  $x_i \sim N_m(\mu_i, \Omega)$ , 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  $M' = (\mu_1, \dots, \mu_n)$ . 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. The following theorem gives some generalizations of the results obtained in Sections 9.4 and 9.5 regarding quadratic forms to these generalized quadratic forms.

**Theorem 9.24.** Let  $X'$  be an  $m \times n$  matrix whose columns are independently distributed, with the  $i$ th column having the  $N_m(\mu_i, \Omega)$  distribution, where  $\Omega$  is positive definite. Suppose that  $A$  and  $B$  are  $n \times n$  symmetric matrices while  $C$  is  $k \times n$ . Let  $M' = (\mu_1, \dots, \mu_n)$ ,  $\Phi = \frac{1}{2}M'AM$ , 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 there exists an  $m \times r$  matrix  $Y'$  such that  $X'AX = Y'Y$ , where the columns of  $Y'$  are independently distributed each having a normal distribution with the same covariance matrix  $\Omega$ , 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  $r$ , there must exist an  $n \times r$  matrix  $P$  satisfying  $A = PP'$  and  $P'P = I_r$ . Consequently,  $X'AX = Y'Y$ , where the  $m \times r$  matrix  $Y' = X'P$  so that

$$\begin{aligned} \text{vec}(Y') &= \text{vec}(X'P) = (P' \otimes I_m)\text{vec}(X') \\ &\sim N_{mr}((P' \otimes I_m)\text{vec}(M'), (P' \otimes I_m)(I_n \otimes \Omega)(P \otimes I_m)) \\ &\sim N_{mr}(\text{vec}(M'P), (I_r \otimes \Omega)) \end{aligned}$$

But this means that the columns of  $Y'$  are independently and normally distributed, each with covariance matrix  $\Omega$ . Further,

$$\frac{1}{2} E(Y')E(Y) = \frac{1}{2} M'PP'M = \frac{1}{2} M'AM = \Phi,$$

and so (a) follows. To prove (b), note that since  $A$  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  $Q'BQ = D$ . Further,  $AB = (0)$  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  $U = QX$ , we find that

$$X'AX = U'CU = \sum_{i=1}^h c_i u_i u_i', \quad X'BX = U'DU = \sum_{i=h+1}^n d_i u_i u_i',$$

where  $u_i$  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).  $\square$

If the columns of the  $m \times n$  matrix  $X'$  are independent and identically distributed as  $N_m(0, \Omega)$  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 very 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  $M = (0)$ . 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 9.25.** Let the columns of the  $m \times n$  matrix  $X' = (x_1, \dots, x_n)$  be independently and identically distributed with  $E(x_i) = 0$ ,  $\text{var}(x_i) = \Omega$ , and  $E(x_i x_i' \otimes x_i x_i') = \Psi$ . 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, since  $\mathbf{x}_1, \dots, \mathbf{x}_n$  are independent, we have

$$\begin{aligned} \text{var}\{\text{vec}(V)\} &= \text{var}\left\{\text{vec}\left(\sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i'\right)\right\} = \text{var}\left\{\sum_{i=1}^n \text{vec}(\mathbf{x}_i \mathbf{x}_i')\right\} \\ &= \sum_{i=1}^n \text{var}\{\text{vec}(\mathbf{x}_i \mathbf{x}_i')\} \\ &= \sum_{i=1}^n \text{var}(\mathbf{x}_i \otimes \mathbf{x}_i) = \sum_{i=1}^n \{E(\mathbf{x}_i \mathbf{x}_i' \otimes \mathbf{x}_i \mathbf{x}_i') - E(\mathbf{x}_i \otimes \mathbf{x}_i)E(\mathbf{x}_i' \otimes \mathbf{x}_i')\} \\ &= \sum_{i=1}^n \{\Psi - \text{vec}(\Omega)\text{vec}(\Omega)'\} = n\{\Psi - \text{vec}(\Omega)\text{vec}(\Omega)'\}. \quad \square \end{aligned}$$

The expression for  $\text{var}\{\text{vec}(V)\}$  simplifies when  $V$  has a Wishart distribution due to the special structure of the fourth moments of the normal distribution. This simplified expression is given in our next theorem. Note that although this theorem is stated for normally distributed columns, the first result given applies to the general case as well.

**Theorem 9.26.** 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  $E(X) = (\mathbf{0})$  and  $E(X'X) = n\Omega$  from the previous theorem, it follows that

$$E(V) = E(X'X + X'M + M'X + M'M) = E(X'X) + M'M = n\Omega + M'M$$

Proceeding as in the proof of Theorem 9.25, we obtain

$$\text{var}\{\text{vec}(V)\} = \sum_{i=1}^n \text{var}\{(x_i + \mu_i) \otimes (x_i + \mu_i)\} \quad (9.13)$$

But

$$\begin{aligned} (x_i + \mu_i) \otimes (x_i + \mu_i) &= x_i \otimes x_i + x_i \otimes \mu_i + \mu_i \otimes x_i + \mu_i \otimes \mu_i \\ &= x_i \otimes x_i + (I_m + K_{mm})(x_i \otimes \mu_i) + \mu_i \otimes \mu_i \\ &= x_i \otimes x_i + 2N_m(I_m \otimes \mu_i)x_i + \mu_i \otimes \mu_i \end{aligned}$$

Since all first and third order moments of  $x_i$  are equal to 0,  $x_i \otimes x_i$  and  $x_i$  are uncorrelated, and so using Theorem 9.20 and Problem 7.52, we find that

$$\begin{aligned} \text{var}\{(x_i + \mu_i) \otimes (x_i + \mu_i)\} &= \text{var}(x_i \otimes x_i) + \text{var}\{2N_m(I_m \otimes \mu_i)x_i\} \\ &= 2N_m(\Omega \otimes \Omega) + 4N_m(I_m \otimes \mu_i)\Omega(I_m \otimes \mu_i')N_m \\ &= 2N_m(\Omega \otimes \Omega) + 4N_m(\Omega \otimes \mu_i \mu_i')N_m \\ &= 2N_m(\Omega \otimes \Omega + \Omega \otimes \mu_i \mu_i' + \mu_i \mu_i' \otimes \Omega) \quad (9.14) \end{aligned}$$

Now substituting (9.14) in (9.13) and simplifying, we obtain (b).  $\square$

**Example 9.8.** In Examples 9.3 and 9.6 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  $x_1, \dots, x_n$  are independently distributed with  $x_i \sim N_m(\mu, \Omega)$  for each  $i$ , and define  $X'$  to be the  $m \times n$  matrix  $(x_1, \dots, x_n)$ . Then the sample mean vector and sample covariance matrix can be expressed as

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} X' \mathbf{1}_n,$$

and

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

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 9.24(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  $W_m(\Omega, n - 1)$ . Further, using Theorem 9.24(c), we see that  $S$  and  $\bar{x}$  are independently distributed since

$$\mathbf{1}'_n (I_n - n^{-1} \mathbf{1}_n \mathbf{1}'_n) = (\mathbf{1}'_n - \mathbf{1}'_n) = \mathbf{0}'$$

In addition, it follows from Theorem 9.26 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 utilizing  $v(S)$ . Since  $v(S) = D_m^+ \text{vec}(S)$ , where  $D_m$  is the duplication matrix discussed in Section 7.8, we have

$$\text{var}\{v(S)\} = \frac{2}{n - 1} D_m^+ N_m(\Omega \otimes \Omega) N_m D_m^{+'}$$

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  $s = (s_{11}, \dots, s_{mm})'$ . Expressions for the mean vector and covariance matrix of  $s$  are easily obtained from the formulas in this example since  $s = w(S) = \Psi_m \text{vec}(S)$  as seen in Problem 7.45, where

$$\Psi_m = \sum_{i=1}^m e_{i,m} (e_{i,m} \otimes e_{i,m})'$$

Thus, using the properties of  $\Psi_m$  obtained in Problem 7.45, we find that

$$\begin{aligned} E(s) &= \Psi_m \text{vec}\{E(S)\} = \Psi_m \text{vec}(\Omega) = w(\Omega), \\ \text{var}(s) &= \Psi_m \text{var}\{\text{vec}(S)\} \Psi_m' = \Psi_m \left\{ \frac{2}{n - 1} N_m(\Omega \otimes \Omega) N_m \right\} \Psi_m' \\ &= \frac{2}{n - 1} \Psi_m (\Omega \otimes \Omega) \Psi_m' = \frac{2}{n - 1} (\Omega \odot \Omega), \end{aligned}$$

where  $\odot$  is the Hadamard product.

**Example 9.9.** The perturbation formulas for eigenvalues and eigenvectors of a symmetric matrix obtained in Section 8.6 can be used to approximate the distributions of an eigenvalue or an eigenvector of a matrix having a Wishart distribution. One important application in statistics that utilizes 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, while 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 9.8, we have, asymptotically,

$$\sqrt{n-1} \{ \text{vec}(S) - \text{vec}(\Omega) \} \sim N_{m^2}(\mathbf{0}, 2N_m(\Omega \otimes \Omega)),$$

where  $\Omega$  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  $\gamma_i$  is a normalized eigenvector of  $S = \Omega + W$  corresponding to the  $i$ th largest eigenvalue  $\lambda_i$ , while  $q_i$  is a normalized eigenvector of  $\Omega$  corresponding to its  $i$ th largest eigenvalue  $x_i$ . Now if  $x_i$  is a distinct eigenvalue of  $\Omega$ , then we have the first-order approximations from Section 8.6

$$\begin{aligned} \lambda_i &= x_i + q_i' W q_i = x_i + (q_i' \otimes q_i') \text{vec}(W), \\ \gamma_i &= q_i - (\Omega - x_i I_m)^+ W q_i = q_i - \{ q_i' \otimes (\Omega - x_i I_m)^+ \} \text{vec}(W) \end{aligned} \quad (9.15)$$

Thus, the asymptotic normality of  $a_i = \sqrt{n-1}(\lambda_i - x_i)$  follows from the asymptotic normality of  $\text{vec}(W_*)$ . Further, we have, asymptotically,

$$\begin{aligned} E(a_i) &= (q_i' \otimes q_i') E\{ \text{vec}(W_*) \} = (q_i' \otimes q_i') \mathbf{0} = \mathbf{0}, \\ \text{var}(a_i) &= (q_i' \otimes q_i') (\text{var}\{ \text{vec}(W_*) \}) (q_i \otimes q_i) \\ &= (q_i' \otimes q_i') (2N_m(\Omega \otimes \Omega)) (q_i \otimes q_i) = 2(q_i' \Omega q_i \otimes q_i' \Omega q_i) = 2x_i^2; \end{aligned}$$

that is, for large  $n$ ,  $\lambda_i \sim N(x_i, 2x_i^2/(n-1))$ , approximately. Similarly,  $b_i = \sqrt{n-1}(\gamma_i - q_i)$  is asymptotically normal with

$$\begin{aligned} E(b_i) &= - \{ q_i' \otimes (\Omega - x_i I_m)^+ \} E\{ \text{vec}(W_*) \} = - \{ q_i' \otimes (\Omega - x_i I_m)^+ \} \mathbf{0} = \mathbf{0}, \\ \Phi = \text{var}(b_i) &= \{ q_i' \otimes (\Omega - x_i I_m)^+ \} \{ \text{var}\{ \text{vec}(W_*) \} \} \{ q_i' \otimes (\Omega - x_i I_m)^+ \}' \\ &= \{ q_i' \otimes (\Omega - x_i I_m)^+ \} \{ 2N_m(\Omega \otimes \Omega) \} \{ q_i \otimes (\Omega - x_i I_m)^+ \} \end{aligned}$$



$$\begin{aligned}
 &= \{(\Omega - x_i I_m)^+ \otimes \mathbf{q}'_i + \mathbf{q}'_i \otimes (\Omega - x_i I_m)^+\} \{(\Omega \otimes \Omega)\} \{\mathbf{q}_i \otimes (\Omega - x_i I_m)^+\} \\
 &= \mathbf{q}'_i \Omega \mathbf{q}_i \otimes (\Omega - x_i I_m)^+ \Omega (\Omega - x_i I_m)^+ = \lambda_i \left\{ \sum_{j \neq i} \frac{x_j}{(x_j - x_i)^2} \mathbf{q}_j \mathbf{q}'_j \right\},
 \end{aligned}$$

and so for large  $n$ ,  $\boldsymbol{\gamma}_i \sim N_m(\mathbf{q}_i, (n-1)^{-1} \Phi)$ , approximately. While the first-order approximations in (9.15) can be used to obtain the asymptotic distributions, higher-order approximations, such as those given in Theorem 8.5, can be used 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, due to the asymptotic normality of  $a_i$ , is asymptotically chi-squared with one degree of freedom. The mean of this chi-squared distribution is 1, while 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. The higher-order approximations of  $\lambda_i$  can be used 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  $t$ . This type of adjustment of asymptotically chi-squared statistics is commonly referred to as a Bartlett adjustment [Bartlett (1937, 1947)]. Some 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 our next example.

**Example 9.10.** A multivariate analysis of variance, such as the multivariate one-way classification model discussed in Example 3.14, utilizes 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 9.30). We will show that if the rank of the noncentrality matrix  $\Phi$  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

$$P\{\lambda_{r+i}(BW^{-1}) > c\} \leq P\{\lambda_i(V_1 V_2^{-1}) > c\},$$

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  $\text{rank}(\Phi) = r$ , there exists an  $r \times m$  matrix  $T$  such that  $\frac{1}{2}T'T = \Phi$ . If we define the  $m \times b$  matrix  $M' = (T' \quad 0)$ , then since  $\frac{1}{2}M'M = \Phi$  and  $B \sim W_m(I_m, b, \Phi)$ , it follows that  $B$  can be expressed as  $B = X'X$ , where  $X'$  is a  $m \times b$  matrix for which  $\text{vec}(X') \sim N_{bm}(\text{vec}(M'), I_b \otimes I_m)$ . Partitioning  $X'$  as  $X' = (X'_1 \quad X'_2)$ , where  $X'_1$  is  $m \times r$ , we find that

$$B = X'_1 X_1 + X'_2 X_2 = B_1 + B_2,$$

where  $B_1 \sim W_m(I_m, r, \Phi)$  and  $B_2 \sim W_m(I_m, b-r)$  since  $\text{vec}(X'_1) \sim N_{rm}(\text{vec}(T'), I_r \otimes I_m)$  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

$$S_1(B_1) = \{B_2, W: \lambda_{r+i}(BW^{-1}) > c\},$$

$$S_2(B_1) = \{B_2, W: \lambda_i\{(F'B_2 F)(F'WF)^{-1}\} > c\}$$

It follows from Problem 3.32(a) that

$$\lambda_i\{(F'BF)(F'WF)^{-1}\} = \lambda_i\{(F'B_2 F)(F'WF)^{-1}\} \geq \lambda_{r+i}(BW^{-1}),$$

so for each fixed  $B_1, S_1(B_1) \subseteq S_2(B_1)$ , and it can be easily verified that  $V_1 = F'B_2F \sim W_{m-r}(I_{m-r}, b-r)$  and  $V_2 = F'WF \sim W_{m-r}(I_{m-r}, w)$ . 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

$$\int_{S_1(B_1)} g(W)f_2(B_2) dW dB_2 \leq \int_{S_2(B_1)} g(W)f_2(B_2) dW dB_2 = P\{\lambda_i(V_1V_2^{-1}) > c\}$$

If we also define the sets

$$C_1 = \{B_1, B_2, W: \lambda_{r+i}(BW^{-1}) > c\}$$

$$C_2 = \{B_1: B_1 \text{ positive definite}\},$$

then the desired result follows since

$$\begin{aligned} P\{\lambda_{r+i}(BW^{-1}) > c\} &= \int_{C_1} g(W)f_1(B_1)f_2(B_2) dW dB_1 dB_2 \\ &= \int_{C_2} \left\{ \int_{S_1(B_1)} g(W)f_2(B_2) dW dB_2 \right\} f_1(B_1) dB_1 \\ &\leq \int_{C_2} P\{\lambda_i(V_1V_2^{-1}) > c\} f_1(B_1) dB_1 \\ &= P\{\lambda_i(V_1V_2^{-1}) > c\} \end{aligned}$$

The relationship between the sample correlation and covariance matrices and the expression for  $\text{var}\{\text{vec}(S)\}$  given in Example 9.8 can be used to obtain an expression for the asymptotic covariance matrix of  $\text{vec}(R)$ . This is the subject of our final example.

**Example 9.11.** As in Example 9.8, let  $x_1, \dots, x_n$  be independently distributed with  $x_i \sim N_m(\mu, \Omega)$ , 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

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

while the population correlation matrix is given by

$$P = D_{\Omega}^{-1/2} \Omega D_{\Omega}^{-1/2}$$

Note that if we define  $y_i = D_{\Omega}^{-1/2} x_i$ , then  $y_1, \dots, y_n$  are independently distributed with  $y_i \sim N_m(D_{\Omega}^{-1/2} \mu, P)$ . If  $S_*$  is the sample covariance matrix computed from the  $y_i$ 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  $y_i$ s is the same as that computed from the  $x_i$ s. If  $A = S_* - P$ , then the first-order approximation for  $R$  is given by (see Problem 8.15)

$$R = P + A - \frac{1}{2} (P D_A + D_A P),$$

and so

$$\begin{aligned} \text{vec}(R) &= \text{vec}(P) + \text{vec}(A) - \frac{1}{2} \{ \text{vec}(P D_A) + \text{vec}(D_A P) \} \\ &= \text{vec}(P) + \text{vec}(A) - \frac{1}{2} \{ (I_m \otimes P) + (P \otimes I_m) \} \text{vec}(D_A) \\ &= \text{vec}(P) + \left( I_{m^2} - \frac{1}{2} \{ (I_m \otimes P) + (P \otimes I_m) \} \Lambda_m \right) \text{vec}(A), \quad (9.16) \end{aligned}$$

where

$$\Lambda_m = \sum_{i=1}^m (E_{ii} \otimes E_{ii})$$

Thus, since

$$\text{var}\{\text{vec}(A)\} = \text{var}\{\text{vec}(S_*)\} = \frac{2}{n-1} N_m(P \otimes P) N_m,$$

we get the first-order approximation

$$\text{var}\{\text{vec}(R)\} = \frac{2}{n-1} H N_m(P \otimes P) N_m H',$$

where the matrix  $H$  is the premultiplier on  $\text{vec}(A)$  in the last expression given in (9.16). Simplification (see Problem 9.33) leads to

$$\text{var}\{\text{vec}(R)\} = \frac{2}{n-1} N_m \Phi N_m, \tag{9.17}$$

where

$$\Phi = \{I_{m^2} - (I_m \otimes P)\Lambda_m\}(P \otimes P)\{I_{m^2} - \Lambda_m(I_m \otimes P)\}$$

Since  $R$  is symmetric and has each diagonal element equal to one, its redundant and nonrandom elements can be eliminated by utilizing  $\tilde{v}(R)$ . Since  $\tilde{v}(R) = \tilde{L}_m \text{vec}(R)$ , where  $\tilde{L}_m$  is the matrix discussed in Section 7.8, we find that the asymptotic covariance matrix of  $\tilde{v}(R)$  is given by

$$\text{var}\{\tilde{v}(R)\} = \frac{2}{n-1} \tilde{L}_m N_m \Phi N_m \tilde{L}_m'$$

The Hadamard product and its associated properties can be useful in analyses involving the manipulation of  $\Phi$  since

$$\begin{aligned} \Phi &= P \otimes P - (I_m \otimes P)\Lambda_m(P \otimes P) - (P \otimes P)\Lambda_m(I_m \otimes P) \\ &\quad + (I_m \otimes P)\Lambda_m(P \otimes P)\Lambda_m(I_m \otimes P), \end{aligned}$$

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**

1. We saw in the proof of Theorem 9.1 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.
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.
3. Let  $A$  be an  $m \times n$  matrix. Show that each of the following matrices is idempotent.

- (a)  $AA^-$ .  
 (b)  $A^-A$ .  
 (c)  $A(A'A)^-A'$ .
4. Determine the class of  $m \times 1$  vectors  $\{x\}$ , for which  $xx'$  is idempotent.
5. Determine constants  $a$ ,  $b$ , and  $c$  so that each of the following is an idempotent matrix.  
 (a)  $a\mathbf{1}_m\mathbf{1}'_m$ .  
 (b)  $bI_m + c\mathbf{1}_m\mathbf{1}'_m$ .
6. Let  $A$  be an  $m \times n$  matrix with  $\text{rank}(A) = m$ . Show that  $A'(AA')^{-1}A$  is symmetric and idempotent and find its rank.
7. Let  $A$  and  $B$  be  $m \times m$  matrices. Show that if  $B$  is nonsingular and  $AB$  is idempotent, then  $BA$  is also idempotent.
8. Show that if  $A$  is an  $m \times m$  matrix and  $A^2 = mA$  for some scalar  $m$ , then

$$\text{tr}(A) = m \text{rank}(A)$$

9. Give an example of a collection of matrices  $A_1, \dots, A_k$  that satisfies conditions (a) and (d) of Corollary 9.7.1, but does not satisfy conditions (b) and (c). Similarly, find a collection of matrices that satisfies conditions (c) and (d) but does not satisfy conditions (a) and (b).
10. Prove Theorem 9.11.
11. Let  $A$  be an  $m \times m$  symmetric matrix with  $r = \text{rank}(A)$  and suppose that  $x \sim N_m(\mathbf{0}, I_m)$ . Show that the distribution of  $x'Ax$  can be expressed as a linear combination of  $r$  independent chi-squared random variables, each with 1 degree of freedom. What are the coefficients in this linear combination when  $A$  is idempotent?
12. Extend the result of Problem 11 to the situation in which  $x \sim N_m(\mathbf{0}, \Omega)$ , where  $\Omega$  is nonnegative definite; that is, show that if  $A$  is a symmetric matrix, then  $x'Ax$  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?
13. Let  $x_1, \dots, x_n$  be a random sample from a normal distribution with mean  $\mu$  and variance  $\sigma^2$ , 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  $t$ ?

14. Suppose that  $\mathbf{x} \sim N_n(\boldsymbol{\mu}, \boldsymbol{\Omega})$ , where  $\boldsymbol{\Omega}$  is positive definite. Partition  $\mathbf{x}$ ,  $\boldsymbol{\mu}$ , and  $\boldsymbol{\Omega}$  as

$$\mathbf{x} = \begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix}, \quad \boldsymbol{\mu} = \begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \quad \boldsymbol{\Omega} = \begin{bmatrix} \Omega_{11} & \Omega_{12} \\ \Omega'_{12} & \Omega_{22} \end{bmatrix},$$

where  $\mathbf{x}_1$  is  $r \times 1$  and  $\mathbf{x}_2$  is  $(n - r) \times 1$ . Show that

- (a)  $t_1 = (\mathbf{x}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Omega}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \sim \chi_r^2$ ,
- (b)  $t_2 = (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Omega}^{-1} (\mathbf{x} - \boldsymbol{\mu}) - (\mathbf{x}_1 - \boldsymbol{\mu}_1)' \boldsymbol{\Omega}_{11}^{-1} (\mathbf{x}_1 - \boldsymbol{\mu}_1) \sim \chi_{n-r}^2$ ,
- (c)  $t_1$  and  $t_2$  are independently distributed.

15. Prove Theorem 9.14.

16. 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  $\mu_i$ s are nonnegative constants satisfying  $\mu_1 + \dots + \mu_m = 1$ . Let  $\mathbf{x} = (x_1, \dots, x_m)'$ ,  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_m)'$ , and  $\boldsymbol{\Omega} = D - \boldsymbol{\mu}\boldsymbol{\mu}'$ , where  $D = \text{diag}(\mu_1, \dots, \mu_m)$ .

- (a) Show that  $\boldsymbol{\Omega}$  is a singular matrix.
  - (b) Show that if  $\sqrt{n}(\mathbf{x} - \boldsymbol{\mu}) \sim N_m(\mathbf{0}, \boldsymbol{\Omega})$ , then  $t \sim \chi_{m-1}^2$ .
17. Suppose that  $\mathbf{x} \sim N_4(\mathbf{0}, \mathbf{I}_4)$  and consider the three functions of the components of  $\mathbf{x}$  given by

$$t_1 = \frac{1}{4} (x_1 + x_2 + x_3 + x_4)^2 + \frac{1}{2} (x_1 - x_2)^2,$$

$$t_2 = \frac{1}{12} (x_1 + x_2 + x_3 - 3x_4)^2,$$

$$t_3 = (x_1 + x_2 - 2x_3)^2 + (x_3 - x_4)^2$$

- (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?

18. Suppose that  $\mathbf{x} \sim N_4(\boldsymbol{\mu}, \boldsymbol{\Omega})$ , where  $\boldsymbol{\mu} = (1, -1, 1, -1)'$  and  $\boldsymbol{\Omega} = \mathbf{I}_4 + \mathbf{1}_4\mathbf{1}_4'$ . Define

$$t_1 = \frac{1}{2} (x_1 - x_2)^2 + \frac{1}{2} (x_3 - x_4)^2,$$

$$t_2 = \frac{1}{2} (x_1 + x_2)^2 + \frac{1}{2} (x_3 + x_4)^2$$

- (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?
19. Prove Theorem 9.15.

20. Prove Theorem 9.16.

21. The purpose of this exercise is to generalize the results of Example 9.5 to a test of the hypothesis that  $H\boldsymbol{\beta} = \mathbf{c}$ , where  $H$  is an  $m_2 \times m$  matrix having rank  $m_2$  and  $\mathbf{c}$  is an  $m_2 \times 1$  vector; Example 9.5 dealt with the special case in which  $H = ((0) \quad \mathbf{I}_{m_2})$  and  $\mathbf{c} = \mathbf{0}$ . 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

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

where  $\mathbf{y}_* = \mathbf{y} - XH'(HH')^{-1}\mathbf{c}$ ,  $X_* = XG'(GG')^{-1}$ , and  $\boldsymbol{\beta}_* = G\boldsymbol{\beta}$ . 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.

22. Suppose that  $\mathbf{x} \sim N_m(\mathbf{0}, \boldsymbol{\Omega})$ , where  $r = \text{rank}(\boldsymbol{\Omega}) < m$ . If  $T$  is any  $m \times r$  matrix satisfying  $TT' = \boldsymbol{\Omega}$ , and  $\mathbf{z} \sim N_r(\mathbf{0}, \mathbf{I}_r)$ , then  $\mathbf{x}$  is distributed the same as  $T\mathbf{z}$ . Use this to show that the formulas given in Theorem 9.21 for positive definite  $\boldsymbol{\Omega}$  also hold when  $\boldsymbol{\Omega}$  is positive semidefinite.

23. Let  $\mathbf{x} \sim N_m(\mathbf{0}, \mathbf{I}_m)$ . Use the fact that the first six moments of the standard normal distribution are 0, 1, 0, 3, 0, and 15 to show that



$$\begin{aligned}
 E(\mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}' \otimes \mathbf{x}\mathbf{x}') &= \mathbf{I}_m^3 + \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m (\mathbf{I}_m \otimes T_{ij} \otimes T_{ij} + T_{ij} \otimes \mathbf{I}_m \otimes T_{ij} \\
 &\quad + T_{ij} \otimes T_{ij} \otimes \mathbf{I}_m) \\
 &\quad + \sum_{i=1}^m \sum_{j=1}^m \sum_{k=1}^m (T_{ij} \otimes T_{ik} \otimes T_{jk}),
 \end{aligned}$$

where  $T_{ij} = E_{ij} + E_{ji}$ .

24. Let  $A$ ,  $B$ , and  $C$  be  $m \times m$  symmetric matrices and suppose that  $\mathbf{x} \sim N_m(\mathbf{0}, \mathbf{I}_m)$ .

(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 the previous exercise to derive the formula for  $E(\mathbf{x}'A\mathbf{x}\mathbf{x}'B\mathbf{x}\mathbf{x}'C\mathbf{x})$  given in Theorem 9.23.

25. Let  $\mathbf{x} \sim N_m(\boldsymbol{\mu}, \boldsymbol{\Omega})$ , where  $\boldsymbol{\Omega}$  is positive definite.

(a) Using Theorem 9.20, show that

$$\text{var}(\mathbf{x} \otimes \mathbf{x}) = 2N_m(\boldsymbol{\Omega} \otimes \boldsymbol{\Omega} + \boldsymbol{\Omega} \otimes \boldsymbol{\mu}\boldsymbol{\mu}' + \boldsymbol{\mu}\boldsymbol{\mu}' \otimes \boldsymbol{\Omega})$$

(b) Show that the matrix  $(\boldsymbol{\Omega} \otimes \boldsymbol{\Omega} + \boldsymbol{\Omega} \otimes \boldsymbol{\mu}\boldsymbol{\mu}' + \boldsymbol{\mu}\boldsymbol{\mu}' \otimes \boldsymbol{\Omega})$  is nonsingular.

(c) Determine the eigenvalues of  $N_m$ . Use these along with part (b) to show that  $\text{rank}\{\text{var}(\mathbf{x} \otimes \mathbf{x})\} = m(m + 1)/2$ .

26. 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)  $E(\mathbf{x} \otimes \mathbf{y} \otimes \mathbf{x} \otimes \mathbf{y}) = (\mathbf{I}_m \otimes K_{nm} \otimes \mathbf{I}_n)\{\text{vec}(V_1) \otimes \text{vec}(V_2)\}$ ,

(e)  $\text{var}(\mathbf{x} \otimes \mathbf{y}) = V_1 \otimes V_2 - \boldsymbol{\mu}_1\boldsymbol{\mu}_1' \otimes \boldsymbol{\mu}_2\boldsymbol{\mu}_2'$ .

27. 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  $\mathbf{x} \sim N_m(\mathbf{0}, \boldsymbol{\Omega})$ , show that

(a)  $E(\mathbf{x}'A\mathbf{a}\mathbf{x}'B\mathbf{b}) = \mathbf{a}'A\boldsymbol{\Omega}B\mathbf{b}$ ,

(b)  $E(\mathbf{x}'A\mathbf{a}\mathbf{x}'B\mathbf{b}\mathbf{x}'C\mathbf{x}) = \mathbf{a}'A\boldsymbol{\Omega}B\mathbf{b} \text{tr}(\boldsymbol{\Omega}C) + 2\mathbf{a}'A\boldsymbol{\Omega}C\boldsymbol{\Omega}B\mathbf{b}$ .

28. Suppose that  $x \sim N_4(\boldsymbol{\mu}, \boldsymbol{\Omega})$ , where  $\boldsymbol{\mu} = \mathbf{1}_4$  and  $\boldsymbol{\Omega} = 4I_4 + \mathbf{1}_4\mathbf{1}'_4$ . 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 9.22 to find

- (a)  $\text{var}(t_1)$ ,  
 (b)  $\text{var}(t_2)$ ,  
 (c)  $\text{cov}(t_1, t_2)$ .

29. Verify the formulas given at the end of Example 9.7 for  $E(t_1)$  and  $\text{var}(t_1)$ .
30. Suppose that the  $m \times 1$  vectors  $\{y_{ij}, 1 \leq i \leq k, 1 \leq j \leq n_i\}$  are independently distributed with  $y_{ij} \sim N_m(\boldsymbol{\mu}_i, \boldsymbol{\Omega})$ . A multivariate analysis of variance utilizes the matrices (Example 3.14)

$$B = \sum_{i=1}^k n_i(\bar{y}_i - \bar{y})(\bar{y}_i - \bar{y})', \quad W = \sum_{i=1}^k \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_i)(y_{ij} - \bar{y}_i)'$$

where

$$\bar{y}_i = \sum_{j=1}^{n_i} \frac{y_{ij}}{n_i}, \quad \bar{y} = \sum_{i=1}^k \frac{n_i \bar{y}_i}{n}, \quad n = \sum_{i=1}^k n_i$$

Use Theorem 9.24 to show that  $W$  and  $B$  are independently distributed,  $W \sim W_m(\boldsymbol{\Omega}, w)$ , and  $B \sim W_m(\boldsymbol{\Omega}, b, \Phi)$ , where  $w = n - k$ ,  $b = k - 1$ , and

$$\Phi = \frac{1}{2} \sum_{i=1}^k n_i(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})', \quad \bar{\boldsymbol{\mu}} = \sum_{i=1}^k \frac{n_i \boldsymbol{\mu}_i}{n}$$

31. Let  $X' = (x_1, \dots, x_n)$  be an  $m \times n$  matrix, where  $x_1, \dots, x_n$  are independent and  $x_i \sim N_m(\mathbf{0}, \boldsymbol{\Omega})$  for each  $i$ . Show that

$$E(X \otimes X \otimes X \otimes X) = \{\text{vec}(I_n) \otimes \text{vec}(I_n)\} \{\text{vec}(\boldsymbol{\Omega}) \otimes \text{vec}(\boldsymbol{\Omega})\}' \\ + \text{vec}(I_n \otimes I_n) \{\text{vec}(\boldsymbol{\Omega} \otimes \boldsymbol{\Omega})\}' + \text{vec}(K_{nn}) \\ \cdot [\text{vec}\{K_{mm}(\boldsymbol{\Omega} \otimes \boldsymbol{\Omega})\}]'$$

32. Suppose that the columns of  $X' = (x_1, \dots, x_n)$  are independently distributed with  $x_i \sim N_m(\mu_i, \Omega)$ . Let  $A$  be an  $m \times m$  symmetric matrix, and let  $M' = (\mu_1, \dots, \mu_n)$ . Use the spectral decomposition of  $A$  to show that

(a)  $E(X'AX) = \text{tr}(A)\Omega + M'AM,$

(b)  $\text{var}\{\text{vec}(X'AX)\} = 2N_m\{\text{tr}(A^2)(\Omega \otimes \Omega) + \Omega \otimes M'A^2M + M'A^2M \otimes \Omega\}$

33. Use the results of Problems 7.45(e) and 7.52 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 (9.17).

# 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.
- Barndorff-Nielsen, O. E. and Cox, D. R. (1994). *Inference 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. (1974). *Generalized Inverses: Theory and Applications*. John Wiley, New York.
- Berman, A. and Plemmons, R. J. (1994). *Nonnegative Matrices in the Mathematical Sciences*. Society for Industrial and Applied Mathematics, Philadelphia.
- Bhattacharya, R. N. and Waymire, E. C. (1990). *Stochastic Processes with Applications*. John Wiley, New York.
- 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. (1990). *Statistical Inference*. Wadsworth & Brooks/Cole, 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. (1979). *Circulant Matrices*. John Wiley, New York.
- Duff, I. S., Erisman, A. M., and Reid, J. K. (1986). *Direct Methods for Sparse Matrices*. Oxford University Press.
- 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.
- 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. (1989). *Matrix Computations*. 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.
- Hageman, L. A. and Young, D. M. (1981). *Applied Iterative Methods*. Academic Press, New York.
- Hammarling, S. J. (1970). *Latent Roots and Latent Vectors*. University of Toronto Press.
- 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.
- Horn, R. A. and Johnson, C. R. (1985). *Matrix Analysis*. Cambridge University Press.
- Horn, R. A. and Johnson, C. R. (1991). *Topics in Matrix Analysis*. Cambridge University Press.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology*, 24, 417–441, 498–520.
- Huberty, C. J. (1994). *Applied Discriminant Analysis*. John Wiley, New York.
- Jackson, J. E. (1991). *A User's Guide to Principal Components*. John Wiley, New York.
- Jolliffe, I. T. (1986). *Principal Component Analysis*. 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.
- Khuri, A. (1993). *Advanced Calculus with Applications in Statistics*. John Wiley, New York.
- Krzanowski, W. J. (1988). *Principles of Multivariate Analysis: A User's Perspective*. Clarendon Press, Oxford.
- 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.

- 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. (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.
- 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.
- 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. (1992). *Discriminant Analysis and Statistical Pattern Recognition*. John Wiley, New York.
- Medhi, J. (1994). *Stochastic Processes*. John Wiley, New York.
- 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. (1990). *Multivariate Statistical Methods*. 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.
- Neter, J., Wasserman, W., and Kutner, M. H. (1985). *Applied Linear Statistical Models: Regression, Analysis of Variance, and Experimental Design*. Irwin, Homewood, IL.
- Olkin, I. and Tomsky, 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). *Solution of Equations in Euclidean and Banach Spaces*. Academic Press, New York.
- 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.
- Poincaré, H. (1890). Sur les équations aux dérivées partielles de la physique mathématique. *American Journal of Mathematics*, **12**, 211–294.
- Press, W. H., Flannery, B. P., Teukolsky, S. A., and Vetterline, W. T. (1992). *Numerical Recipes in FORTRAN: The Art of Scientific Computing*. Cambridge University Press.
- 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. and Mitra, S. K. (1971). *Generalized Inverse of Matrices and Its Applications*, John Wiley, New York.
- Rockafellar, R. T. (1970). *Convex Analysis*. Princeton University Press.
- 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.
- 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. (1973). *Non-negative Matrices: An Introduction to Theory and Applications*. John Wiley, New York.
- Srivastava, M. S. and Khatri, C. G. (1979). *An Introduction to Multivariate Analysis*. North-Holland, New York.
- 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 the latent vector of the Wishart and multivariate F matrices. *Journal of Multivariate Analysis*, **6**, 500–525.
- Taylor, H. M. and Karlin, S. (1984). *An Introduction to Stochastic Modeling*. Academic Press, Orlando.
- Young, D. M. (1971). *Iterative Solution of Large Linear Systems*. Academic Press, New York.





# Index

- Accumulation point, 70  
Adjoint, 8  
Analysis of variance, 120. *See also*  
    One-way classification model; Two-way  
    classification model
- Bartlett adjustment, 406  
Basis, 41–43  
    orthonormal, 48–52  
Bilinear form, 15  
Block diagonal matrix, 12  
Boundary point, 72
- Canonical variate analysis, 107, 154–155,  
    406–407  
Cauchy–Schwarz inequality, 35  
Cayley–Hamilton theorem, 93  
Chain rule, 324, 327  
Characteristic equation, 85  
Characteristic root, 84. *See* Eigenvalue  
Characteristic vector, 84. *See* Eigenvector  
Chi-squared distribution:  
    central, 20–21  
    and Moore–Penrose inverse, 179–180  
    noncentral, 21  
    and quadratic forms, 378–384  
Cholesky decomposition, 139  
Circulant matrix, 300–304  
Closure, 70  
Cochran’s theorem, 374–378  
Cofactor, 5, 8  
    expansion formula for determinant, 5–6  
Column space, 43  
Commutation matrix, 276–283  
    eigenvalues, 281  
    eigenvectors, 317  
Complex matrix, 16–18
- Concave function, *see* Convex function  
Consistent equations, 210–213  
Consistent estimator, 189–190  
Continuity:  
    of determinant, 188  
    of eigenvalues, 103  
    of inverse matrix, 188  
    of Moore–Penrose inverse, 189  
Convex combination, 70  
Convex function, 349–353  
    absolute maximum, 352  
Convex hull, 70  
Convex set, 70–74  
Correlation coefficient, 24  
    maximum squared, 368  
Correlation matrix, 24  
    nonnegative definite, 24  
    sample, 25  
Courant–Fischer min–max theorem, 108–110  
Covariance, 22–23  
    of quadratic forms, 391, 394  
Covariance matrix, 23  
    nonnegative definite, 23  
    sample, 25
- Decomposition:  
    Cholesky, 139  
    Jordan, 147–149  
    LU, 169  
    QR, 140  
    Schur, 149–153  
    singular value, 131–138  
    spectral, 95, 98, 138  
Density function, 19  
Derivative, 323, 325  
    of determinant, 332, 336  
    of eigenvalue, 343

- Derivative (*Continued*)  
 of eigenvector, 343  
 of inverse, 333, 336–337  
 of Moore–Penrose inverse, 333–334, 336–337  
 partial, 325  
 of patterned matrices, 335–337  
 second-order partial, 326  
 of trace, 332  
 of vector function, 327
- Determinant, 5–8  
 continuity of, 188  
 derivative of, 332, 336  
 and eigenvalues, 90  
 expansion formula for, 5–6  
 of partitioned matrix, 249–250
- Diagonalization, 92, 144–147  
 simultaneous, 118, 154–157
- Diagonal matrix, 2
- Differential, 324, 325  
 of determinant, 332  
 of eigenvalue, 343  
 of eigenvector, 343  
 of inverse, 333, 336–337  
 of matrix function, 328  
 of Moore–Penrose inverse, 334–335, 336–337  
 second, 326  
 of trace, 332  
 of vector function, 327
- Dimension of vector space, 41
- Direct sum of matrices, 260–261
- Discriminant analysis, 37
- Distance function, 36  
 Euclidean, 36, 50, 62–63, 141  
 Mahalanobis, 37, 63, 141
- Distance in the metric of, 37
- Duplication matrix, 238–285
- Eigenprojection, 98  
 continuity of, 103
- Eigenspace, 87, 146
- Eigenvalue, 84  
 asymptotic distribution of, 404–406  
 continuity of, 103  
 derivative of, 343  
 distinct, 86  
 extremal properties, 104–110  
 of idempotent matrix, 370–371  
 in the metric of, 118  
 monotonicity, 115  
 multiple, 86  
 of orthogonal matrix, 88  
 perturbation of, 339–343  
 of positive definite matrix, 112  
 of positive semidefinite matrix, 112  
 and rank, 92, 99, 146–147, 153  
 simple, 86  
 of symmetric matrix, 93–102  
 of transpose product, 114–115  
 of triangular matrix, 88
- Eigenvectors, 84  
 asymptotic distribution of, 404–406  
 common, 128, 157  
 derivative of, 343  
 linear independence of, 91  
 of symmetric matrix, 94–96
- Elementary transformations, 13
- Elimination matrices, 285–288
- Estimable function, 230
- Euclidean norm, 36, 37, 158
- Euclidean space, 36
- Euler's formula, 17
- Expected value, 19  
 of quadratic form, 390–398
- F distribution, 21–22
- Fourier matrix, 303–304
- Gauss–Seidel method, 236
- Generalized inverse, 190–196. *See also*  
 Moore–Penrose inverse  
 computation of, 200–203  
 properties, 193
- Gradient, 237
- Gram–Schmidt orthonormalization, 48, 54–55
- Hadamard inequality, 270
- Hadamard matrix, 305–307  
 normalized, 306
- Hadamard product, 266–276  
 eigenvalues of, 274–276  
 as a Kronecker product, 267  
 rank of, 267
- Hermite form, 200
- Hermitian matrix, 18
- Hessian matrix, 326
- Homogeneous system of equations, 219–221
- Hyperplane, 71
- Idempotent matrix, 3, 58–59, 370–374  
 eigenvalues of, 370–371  
 rank of, 370–371  
 symmetric, 372, 373–374  
 trace of, 370–371

- Identity matrix, 2  
 Indefinite matrix, 16  
 Independence (linear), 38–40  
 Independence (stochastic):  
   of quadratic forms, 384–390  
   of random variables, 22  
 Inner product, 34–35  
   Euclidean, 35  
 Interior point, 72  
 Intersection of vector spaces, 67  
 Inverse matrix, 8–11  
   and cofactors, 8–9  
   continuity of, 188  
   derivative of, 333, 336–337  
   of partitioned matrix, 347  
   of a sum, 9–10  
 Irreducible matrix, 294–295
- Jacobian matrix, 327  
 Jacobi method, 236  
 Jensen's inequality, 352–353  
 Jordan decomposition, 147–149
- Kronecker product, 253  
   determinant of, 256  
   eigenvalues of, 255  
   eigenvectors of, 312  
   inverse of, 255  
   Moore–Penrose inverse of, 255  
   rank of, 257  
   trace of, 255
- Lagrange function, 354  
 Lagrange multipliers, 354  
 Lanczos vectors, 238  
 Latent root, 84. *See* Eigenvalue  
 Latent vector, 84. *See* Eigenvector  
 Least squares, *see also* Regression  
   and best linear unbiased estimator,  
     113–114  
   generalized, 141–142, 245  
   in less than full rank models, 58,  
     228–232  
   and multicollinearity, 96–98, 136  
   in multiple regression, 55–58  
   in one-way classification model, 79–80  
   ordinary, 26–28  
   restricted, 80–81, 245  
   in ridge regression, 123  
   in simple linear regression, 50–51  
   and solutions to a system of equations,  
     222–228, 345–346  
   with standardized explanatory variables,  
     64–65  
   weighted, 65–66  
 Least squares inverse, 196–197  
   computation of, 203–204  
 Limit point, 70  
 Linear combination, 33  
 Linear dependence, 38–40  
 Linear equations, 66–67  
   consistency of, 210–213  
   homogeneous system of, 219–221  
   least squares solutions of, 222–228  
   linearly independent solutions to, 217  
   and singular value decomposition,  
     233–235  
   solutions to, 213–219  
   sparse systems of, 235–241  
     direct methods, 235–236  
     iterative methods, 236–241  
   unique solution to, 216  
 Linear independence, 38–40  
 Linear model, 27  
 Linear space, 33  
 Linear transformation, 60–67  
 LU factorization, 169
- Mahalanobis distance, 37, 63, 141  
 Markov chain, 298–300  
 Matrix:  
   block diagonal, 12  
   circulant, 300–304  
   commutation, 276–283  
   complex, 16–18  
   correlation, 24  
   covariance, 23  
   diagonal, 2  
   duplication, 283–285  
   eigenprojection, 98  
   elimination, 285–288  
   Fourier, 303–304  
   Hadamard, 305–307  
   Hermitian, 18  
   Hessian, 326  
   idempotent, 3, 58–59, 370–374  
   identity, 2  
   indefinite, 16  
   irreducible, 294–295  
   Jacobian, 327  
   negative definite, 16  
   negative semidefinite, 16  
   nilpotent, 127, 166  
   nonnegative, 288  
   nonnegative definite, 16

- Matrix (*Continued*)
- nonsingular, 8
  - null, 2
  - order of, 1
  - orthogonal, 14–15
  - partitioned, 11–13
  - permutation, 15
  - positive, 288
  - positive definite, 15–16
  - positive semidefinite, 15–16
  - primitive, 298
  - projection, 52–59
  - reducible, 294–295
  - similar, 144
  - skew-symmetric, 4
  - square root, 16
  - symmetric, 4
  - Toeplitz, 304–305
  - transpose, 3
  - triangular, 2
  - unitary, 18, 150
  - Vandermonde, 307–309
- Matrix function, 327
- Matrix norm, 158
- Euclidean, 158
  - maximum column sum, 158
  - maximum row sum, 158
  - spectral, 158
- Maximum:
- absolute, 344
  - of a concave function, 351
  - conditions for local maximum, 345
  - with equality constraints, 353–360
  - local, 344
- Maximum likelihood estimation, 347–349
- Mean, 19
- sample, 25
- Mean squared error, 163–164
- Mean vector, 22
- differences in, 106–107, 116–117, 154
  - sample, 25
- Minimum, *see* Maximum
- Minor, 5, 13
- leading principal, 311
- Modulus of a complex number, 17
- Moment generating function, 20
- Moments, 19–20
- Moore–Penrose inverse, 171
- of block diagonal matrix, 186
  - computation of, 175, 197–199
  - continuity of, 188–190
  - derivative of, 333–334, 336–337
  - of diagonal matrix, 177
  - existence of, 171–172
  - of a matrix product, 180–185
  - of partitioned matrices, 185–186
  - and projection matrices, 172–173
  - properties, 174–180
  - and quadratic form in normal random vectors, 179–180
  - and rank, 175
  - and singular value decomposition, 172
  - and spectral decomposition, 176–177
  - of a sum, 186–187
  - of a symmetric matrix, 176–178
  - uniqueness of, 171–172
- Multicollinearity, 96–98, 136
- Multinomial distribution, 368
- Multiplicity of an eigenvalue, 86
- Multivariate normal distribution, 25–26, 331–332, 347–349
- Negative definite matrix, 16
- Negative semidefinite matrix, 16
- Nilpotent matrix, 127
- Nonnegative definite matrix, 16
- correlation matrix, 24
  - covariance matrix, 23
- Nonnegative matrix, 288
- irreducible, 294–295
  - eigenvalues of, 296–298
  - eigenvectors of, 296–297
  - primitive, 298
  - spectral radius of, 288
- Nonsingular matrix, 8
- Norm:
- matrix, 157–162
  - vector, 35, 37–38
- Normal distribution:
- multivariate, 25–26, 331–332, 347–349
  - singular, 26, 379
  - univariate, 20
- Normalized vector, 14
- Null matrix, 2
- Null space, 60–61
- Null vector, 2
- One-way classification model:
- multivariate, 119–122, 154, 406–407
  - univariate, 79–80, 119, 228–229, 231–232, 257–258, 385–387, 396–397
- Order:
- of a minor, 13
  - of a square matrix, 1
- Orthogonal complement, 52
- dimension of, 52

- Orthogonal matrix, 14–15  
 Orthogonal vectors, 14  
 Orthonormal basis, 48–52  
 Orthonormal vectors, 14
- Partitioned matrix, 11–13  
   determinant of, 249–250  
   inverse of, 247  
   rank, 46–47
- Pearson's chi-squared statistic, 411
- Permutation matrix, 15
- Perturbation methods, 337–344  
   eigenprojection, 343–344  
   eigenvalue, 339–343  
   matrix inverse, 338–339
- Poincaré separation theorem, 111
- Polar coordinates, 17
- Positive definite matrix, 15–16
- Positive matrix, 288  
   eigenvalues, 289–294  
   eigenvectors, 289–293  
   spectral radius, 288
- Positive semidefinite matrix, 15–16
- Primitive matrix, 298
- Principal components analysis, 107–108, 404
- Probability function, 18
- Projection matrix, 52–59
- Projection (orthogonal), 50
- Quadratic form, 15–16  
   distribution of, 378–384  
   expected value of, 390–398  
   generalized, 399  
   independence of, 384–390  
   and Moore–Penrose inverse, 179–180
- QR factorization, 140
- Random variable, 18–22
- Random vector, 22–26
- Range, 43
- Rank, 13–14  
   and linear independence, 43–47
- Rayleigh quotient, 104
- Reducible matrix, 294–295
- Regression, 26–28. *See also* Least squares  
   best quadratic unbiased estimator, 358–360  
   F test, 387–388  
   generalized least squares, 141–142  
   multiple, 55–58, 248–249  
   principal components, 96–98, 136–138,  
     163  
   ridge, 123  
   simple linear, 50–51  
   with standardized explanatory variables,  
     64–65  
   weighted least squares, 65–66
- Row space, 43
- Saddle point, 345
- Sample correlation matrix, 25  
   asymptotic covariance matrix of, 407–409
- Sample covariance matrix, 25  
   distribution of, 402–403
- Sample mean, 25
- Sample mean vector, 25
- Sample variance, 25  
   distribution, 383–384  
   independent of sample mean, 388–389
- Schur decomposition, 149–153
- Separating hyperplane theorem, 73–74
- Similar matrices, 144
- Simultaneous confidence intervals, 121–122
- Simultaneous diagonalization, 118, 154–157
- Singular value decomposition, 131–138  
   and system of equations, 233–235
- Singular values, 133  
   and eigenvalues, 135
- Skew-symmetric matrix, 4
- Spanning set, 33
- Spectral decomposition, 95, 98, 138
- Spectral radius, 159
- Spectral set, 98
- Square root of a matrix, 16, 138–139
- Stationary point, 345
- Submatrix, 11–13  
   principal, 112
- Subspace, 32
- Sum of squares:  
   for error, 27  
   for treatment, 120
- Sum of vector spaces, 68
- Supporting hyperplane theorem, 73
- Symmetric matrix, 4
- Taylor formula:  
   first-order, 323, 325  
   kth-order, 324, 325  
   vector function, 326
- Toeplitz matrix, 304–305
- Trace, 4–5  
   derivative of, 332  
   and eigenvalues, 90
- Transition probabilities, 299
- Transpose, 3
- Transpose product, 114–115, 142–144
- Triangle inequality, 18, 36

- Triangular matrix, 2
- Two-way classification model, 244, 259–260, 313
- Union–intersection procedure, 121, 367
- Unitary matrix, 18, 150
- Unit vector, 14
- Vandermonde matrix, 307–309
- Variance, 19–20
  - of quadratic form, 391, 394
  - sample, 25
- Vec operator, 261–265
- Vector, 2
  - normalized, 14
  - null, 2
  - orthogonal, 14
  - orthonormal, 14
  - unit, 14
- Vector norm, 35
  - Euclidean norm, 36, 37
  - infinity norm, 37
  - max norm, 37
  - sum norm, 37
- Vector space, 32
  - basis of, 41–43
  - dimension of, 41
  - direct sum, 69
  - Euclidean, 36
  - intersection, 67
  - projection matrix of, 52–59
  - sum, 68
- Wishart distribution, 398–409
  - covariance matrix of, 401
  - mean of, 401
  - and sample covariance matrix, 402–403

## WILEY SERIES IN PROBABILITY AND STATISTICS

ESTABLISHED BY WALTER A. SHEWHART AND SAMUEL S. WILKS

### Editors

*Vic Barnett, Ralph A. Bradley, Nicholas I. Fisher, J. Stuart Hunter,  
J. B. Kadane, David G. Kendall, David W. Scott, Adrian F. M. Smith,  
Jozef L. Teugels, Geoffrey S. Watson*

### *Probability and Statistics*

- ANDERSON · An Introduction to Multivariate Statistical Analysis, *Second Edition*  
\*ANDERSON · The Statistical Analysis of Time Series  
ARNOLD, BALAKRISHNAN, and NAGARAJA · A First Course in Order Statistics  
BACCELLI, COHEN, OLSDER, and QUADRAT · Synchronization and Linearity:  
An Algebra for Discrete Event Systems  
BARTOSZYNSKI and NIEWIADOMSKA-BUGAJ · Probability and Statistical Inference  
BERNARDO and SMITH · Bayesian Statistical Concepts and Theory  
BHATTACHARYYA and JOHNSON · Statistical Concepts and Methods  
BILLINGSLEY · Convergence of Probability Measures  
BILLINGSLEY · Probability and Measure, *Second Edition*  
BOROVKOV · Asymptotic Methods in Queuing Theory  
BRANDT, FRANKEN, and LISEK · Stationary Stochastic Models  
CAINES · Linear Stochastic Systems  
CAIROLI and DALANG · Sequential Stochastic Optimization  
CHEN · Recursive Estimation and Control for Stochastic Systems  
CONSTANTINE · Combinatorial Theory and Statistical Design  
COOK and WEISBERG · An Introduction to Regression Graphics  
COVER and THOMAS · Elements of Information Theory  
CSÖRGÖ and HORVÁTH · Weighted Approximations in Probability Statistics  
\*DOOB · Stochastic Processes  
DUDEWICZ and MISHRA · Modern Mathematical Statistics  
DUPUIS · A Weak Convergence Approach to the Theory of Large Deviations  
ETHIER and KURTZ · Markov Processes: Characterization and Convergence  
FELLER · An Introduction to Probability Theory and Its Applications, Volume I,  
*Third Edition, Revised; Volume II, Second Edition*  
FREEMAN and SMITH · Aspects of Uncertainty: A Tribute to D. V. Lindley  
FULLER · Introduction to Statistical Time Series, *Second Edition*  
FULLER · Measurement Error Models  
GHOSH · Sequential Estimation  
GIFI · Nonlinear Multivariate Analysis  
GUTTORP · Statistical Inference for Branching Processes  
HALD · A History of Probability and Statistics and Their Applications before 1750  
HALL · Introduction to the Theory of Coverage Processes  
HANNAN and DEISTLER · The Statistical Theory of Linear Systems  
HEDAYAT and SINHA · Design and Inference in Finite Population Sampling  
HOEL · Introduction to Mathematical Statistics, *Fifth Edition*  
HUBER · Robust Statistics  
IMAN and CONOVER · A Modern Approach to Statistics  
JUREK and MASON · Operator-Limit Distributions in Probability Theory  
KAUFMAN and ROUSSEEUW · Finding Groups in Data: An Introduction to Cluster  
Analysis  
KOTZ · Leading Personalities in Statistical Sciences from the Seventeenth Century to the  
Present  
LAMPERTI · Probability: A Survey of the Mathematical Theory, *Second Edition*

\*Now available in a lower priced paperback edition in the Wiley Classics Library.

- LARSON · Introduction to Probability Theory and Statistical Inference, *Third Edition*  
LESSLER and KALSBECK · Nonsampling Error in Surveys  
LINDVALL · Lectures on the Coupling Method  
MANTON, WOODBURY, and TOLLEY · Statistical Applications Using Fuzzy Sets  
MARDIA · The Art of Statistical Science: A Tribute to G. S. Watson  
MORGENTHAUER and TUKEY · Configural Polysampling: A Route to Practical Robustness  
MUIRHEAD · Aspects of Multivariate Statistical Theory  
OLIVER and SMITH · Influence Diagrams, Belief Nets and Decision Analysis  
\*PARZEN · Modern Probability Theory and Its Applications  
PRESS · Bayesian Statistics: Principles, Models, and Applications  
PUKELSHEIM · Optimal Experimental Design  
PURI and SEN · Nonparametric Methods in General Linear Models  
PURI, VILAPLANA, and WERTZ · New Perspectives in Theoretical and Applied Statistics  
RAO · Asymptotic Theory of Statistical Inference  
RAO · Linear Statistical Inference and Its Applications, *Second Edition*  
\*RAO and SHANBHAG · Choquet-Deny Type Functional Equations with Applications to Stochastic Models  
RENCHER · Methods of Multivariate Analysis  
ROBERTSON, WRIGHT, and DYKSTRA · Order Restricted Statistical Inference  
ROGERS and WILLIAMS · Diffusions, Markov Processes, and Martingales, Volume I: Foundations, *Second Edition*; Volume II: Itô Calculus  
ROHATGI · An Introduction to Probability Theory and Mathematical Statistics  
ROSS · Stochastic Processes  
RUBINSTEIN · Simulation and the Monte Carlo Method  
RUBINSTEIN and SHAPIRO · Discrete Event Systems: Sensitivity Analysis and Stochastic Optimization by the Score Function Method  
RUZSA and SZEKELY · Algebraic Probability Theory  
SCHEFFE · The Analysis of Variance  
SEBER · Linear Regression Analysis  
SEBER · Multivariate Observations  
SEBER and WILD · Nonlinear Regression  
SERFLING · Approximation Theorems of Mathematical Statistics  
SHORACK and WELLNER · Empirical Processes with Applications to Statistics  
SMALL and MCLEISH · Hilbert Space Methods in Probability and Statistical Inference  
STAPLETON · Linear Statistical Models  
STAUDTE and SHEATHER · Robust Estimation and Testing  
STOYANOV · Counterexamples in Probability  
STYAN · The Collected Papers of T. W. Anderson: 1943–1985  
TANAKA · Time Series Analysis: Nonstationary and Noninvertible Distribution Theory  
THOMPSON and SEBER · Adaptive Sampling  
WELSH · Aspects of Statistical Inference  
WHITTAKER · Graphical Models in Applied Multivariate Statistics  
YANG · The Construction Theory of Denumerable Markov Processes

*Applied Probability and Statistics*

- ABRAHAM and LEDOLTER · Statistical Methods for Forecasting  
AGRESTI · Analysis of Ordinal Categorical Data  
AGRESTI · Categorical Data Analysis  
AGRESTI · An Introduction to Categorical Data Analysis  
ANDERSON and LOYNES · The Teaching of Practical Statistics

\*Now available in a lower priced paperback edition in the Wiley Classics Library.



*Applied Probability and Statistics (Continued)*

- ANDERSON, AUQUIER, HAUCK, OAKES, VANDAELE, and WEISBERG ·  
Statistical Methods for Comparative Studies
- ARMITAGE and DAVID (editors) · Advances in Biometry
- \*ARTHANARI and DODGE · Mathematical Programming in Statistics
- ASMUSSEN · Applied Probability and Queues
- \*BAILEY · The Elements of Stochastic Processes with Applications to the Natural  
Sciences
- BARNETT and LEWIS · Outliers in Statistical Data, *Second Edition*
- BARTHOLOMEW, FORBES, and McLEAN · Statistical Techniques for Manpower  
Planning, *Second Edition*
- 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
- BELSLEY · Conditioning Diagnostics: Collinearity and Weak Data in Regression
- BELSLEY, KUH, and WELSCH · Regression Diagnostics: Identifying Influential  
Data and Sources of Collinearity
- BERRY · Bayesian Analysis in Statistics and Econometrics: Essays in Honor of Arnold  
Zellner
- BERRY, CHALONER, and GEWEKE · Bayesian Analysis in Statistics and  
Econometrics: Essays in Honor of Arnold Zellner
- BHAT · Elements of Applied Stochastic Processes, *Second Edition*
- BHATTACHARYA and WAYMIRE · Stochastic Processes with Applications
- BIEMER, GROVES, LYBERG, MATHIOWETZ, and SUDMAN · Measurement  
Errors in Surveys
- BIRKES and DODGE · Alternative Methods of Regression
- BLOOMFIELD · Fourier Analysis of Time Series: An Introduction
- BOLLEN · Structural Equations with Latent Variables
- BOULEAU · Numerical Methods for Stochastic Processes
- BOX · R. A. Fisher, the Life of a Scientist
- BOX and DRAPER · Empirical Model-Building and Response Surfaces
- BOX and DRAPER · Evolutionary Operation: A Statistical Method for Process  
Improvement
- BOX, HUNTER, and HUNTER · Statistics for Experimenters: An Introduction to  
Design, Data Analysis, and Model Building
- BROWN and HOLLANDER · Statistics: A Biomedical Introduction
- BUCKLEW · Large Deviation Techniques in Decision, Simulation, and Estimation
- BUNKE and BUNKE · Nonlinear Regression, Functional Relations and Robust  
Methods: Statistical Methods of Model Building
- CHATTERJEE and HADI · Sensitivity Analysis in Linear Regression
- CHATTERJEE and PRICE · Regression Analysis by Example, *Second Edition*
- CLARKE and DISNEY · Probability and Random Processes: A First Course with  
Applications, *Second Edition*
- COCHRAN · Sampling Techniques, *Third Edition*
- \*COCHRAN and COX · Experimental Designs, *Second Edition*
- CONOVER · Practical Nonparametric Statistics, *Second Edition*
- CONOVER and IMAN · Introduction to Modern Business Statistics
- CORNELL · Experiments with Mixtures, Designs, Models, and the Analysis of Mixture  
Data, *Second Edition*
- COX · A Handbook of Introductory Statistical Methods
- \*COX · Planning of Experiments
- COX, BINDER, CHINNAPPA, CHRISTIANSON, COLLEDGE, and KOTT ·  
Business Survey Methods
- CRESSIE · Statistics for Spatial Data, *Revised Edition*
- DANIEL · Applications of Statistics to Industrial Experimentation

\*Now available in a lower priced paperback edition in the Wiley Classics Library.

*Applied Probability and Statistics (Continued)*

DANIEL · Biostatistics: A Foundation for Analysis in the Health Sciences, *Sixth Edition*

DAVID · Order Statistics, *Second Edition*

\*DEGROOT, FIENBERG, and KADANE · Statistics and the Law

\*DEMING · Sample Design in Business Research

DILLON and GOLDSTEIN · Multivariate Analysis: Methods and Applications

DODGE and ROMIG · Sampling Inspection Tables, *Second Edition*

DOWDY and WEARDEN · Statistics for Research, *Second Edition*

DRAPER and SMITH · Applied Regression Analysis, *Second Edition*

DUNN · Basic Statistics: A Primer for the Biomedical Sciences, *Second Edition*

DUNN and CLARK · Applied Statistics: Analysis of Variance and Regression, *Second Edition*

ELANDT-JOHNSON and JOHNSON · Survival Models and Data Analysis

EVANS, PEACOCK, and HASTINGS · Statistical Distributions, *Second Edition*

FISHER and VAN BELLE · Biostatistics: A Methodology for the Health Sciences

FLEISS · The Design and Analysis of Clinical Experiments

FLEISS · Statistical Methods for Rates and Proportions, *Second Edition*

FLEMING and HARRINGTON · Counting Processes and Survival Analysis

FLURY · Common Principal Components and Related Multivariate Models

GALLANT · Nonlinear Statistical Models

GLASSERMAN and YAO · Monotone Structure in Discrete-Event Systems

GNANADESIKAN · Methods for Statistical Data Analysis of Multivariate Observations, *Second Edition*

GREENWOOD and NIKULIN · A Guide to Chi-Squared Testing

GROSS and HARRIS · Fundamentals of Queueing Theory, *Second Edition*

GROVES · Survey Errors and Survey Costs

GROVES, BIEMER, LYBERG, MASSEY, NICHOLLS, and WAKSBERG · Telephone Survey Methodology

HAHN and MEEKER · Statistical Intervals: A Guide for Practitioners

HAND · Discrimination and Classification

\*HANSEN, HURWITZ, and MADOW · Sample Survey Methods and Theory, Volume I: Methods and Applications

\*HANSEN, HURWITZ, and MADOW · Sample Survey Methods and Theory, Volume II: Theory

HEIBERGER · Computation for the Analysis of Designed Experiments

HELLER · MACSYMA for Statisticians

HINKELMAN and KEMPTHORNE · Design and Analysis of Experiments, Volume 1: Introduction to Experimental Design

HOAGLIN, MOSTELLER, and TUKEY · Exploratory Approach to 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 Variables

HOEL · Elementary Statistics, *Fifth Edition*

HOGG and KLUGMAN · Loss Distributions

HOLLANDER and WOLFE · Nonparametric Statistical Methods

HOSMER and LEMESHOW · Applied Logistic Regression

HOYLAND and RAUSAND · System Reliability Theory: Models and Statistical Methods

HUBERTY · Applied Discriminant Analysis

IMAN and CONOVER · Modern Business Statistics

JACKSON · A User's Guide to Principle Components

JOHN · Statistical Methods in Engineering and Quality Assurance

\*Now available in a lower priced paperback edition in the Wiley Classics Library.

*Applied Probability and Statistics (Continued)*

JOHNSON · Multivariate Statistical Simulation

JOHNSON and KOTZ · Distributions in Statistics  
Continuous Univariate Distributions—2  
Continuous Multivariate Distributions

JOHNSON, KOTZ, and BALAKRISHNAN · Continuous Univariate Distributions.  
Volume 1, *Second Edition*

JOHNSON, KOTZ, and BALAKRISHNAN · Discrete Multivariate Distributions

JOHNSON, KOTZ, and KEMP · Univariate Discrete Distributions, *Second Edition*

JUDGE, GRIFFITHS, HILL, LÜTKEPOHL, and LEE · The Theory and Practice of  
Econometrics, *Second Edition*

JUDGE, HILL, GRIFFITHS, LÜTKEPOHL, and LEE · Introduction to the Theory and  
Practice of Econometrics, *Second Edition*

JUREČKOVÁ and SEN · Robust Statistical Procedures: Aymptotics and Interrelations

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

KASPRZYK, DUNCAN, KALTON, and SINGH · Panel Surveys

KISH · Statistical Design for Research

\*KISH · Survey Sampling

LAD · Operational Subjective Statistical Methods: A Mathematical, Philosophical, and  
Historical Introduction

LANGE, RYAN, BILLARD, BRILLINGER, CONQUEST, and GREENHOUSE ·  
Case Studies in Biometry

LAWLESS · Statistical Models and Methods for Lifetime Data

LEBART, MORINEAU., and WARWICK · Multivariate Descriptive Statistical  
Analysis: Correspondence Analysis and Related Techniques for Large Matrices

LEE · Statistical Methods for Survival Data Analysis, *Second Edition*

LEPAGE and BILLARD · Exploring the Limits of Bootstrap

LEVY and LEMESHOW · Sampling of Populations: Methods and Applications

LINHART and ZUCCHINI · Model Selection

LITTLE and RUBIN · Statistical Analysis with Missing Data

LYBERG · Survey Measurement

MAGNUS and NEUDECKER · Matrix Differential Calculus with Applications in  
Statistics and Econometrics

MAINDONALD · Statistical Computation

MALLOWS · Design, Data, and Analysis by Some Friends of Cuthbert Daniel

MANN, SCHAFER, and SINGPURWALLA · Methods for Statistical Analysis of  
Reliability and Life Data

MASON, GUNST, and HESS · Statistical Design and Analysis of Experiments with  
Applications to Engineering and Science

McLACHLAN and KRISHNAN · The EM Algorithm and Extensions

McLACHLAN · Discriminant Analysis and Statistical Pattern Recognition

McNEIL · Epidemiological Research Methods

MILLER · Survival Analysis

MONTGOMERY and MYERS · Response Surface Methodology: Process and Product  
in Optimization Using Designed Experiments

MONTGOMERY and PECK · Introduction to Linear Regression Analysis, *Second Edition*

NELSON · Accelerated Testing, Statistical Models, Test Plans, and Data Analyses

NELSON · Applied Life Data Analysis

OCHI · Applied Probability and Stochastic Processes in Engineering and Physical  
Sciences

OKABE, BOOTS, and SUGIHARA · Spatial Tesselations: Concepts and Applications  
of Voronoi Diagrams

\*Now available in a lower priced paperback edition in the Wiley Classics Library.

- OSBORNE · Finite Algorithms in Optimization and Data Analysis  
PANKRATZ · Forecasting with Dynamic Regression Models  
PANKRATZ · Forecasting with Univariate Box-Jenkins Models: Concepts and Cases  
PORT · Theoretical Probability for Applications  
PUTERMAN · Markov Decision Processes: Discrete Stochastic Dynamic Programming  
RACHEV · Probability Metrics and the Stability of Stochastic Models  
RÉNYI · A Diary on Information Theory  
RIPLEY · Spatial Statistics  
RIPLEY · Stochastic Simulation  
ROSS · Introduction to Probability and Statistics for Engineers and Scientists  
ROUSSEEUW and LEROY · Robust Regression and Outlier Detection  
RUBIN · Multiple Imputation for Nonresponse in Surveys  
RYAN · Modern Regression Methods  
RYAN · Statistical Methods for Quality Improvement  
SCHOTT · Matrix Analysis for Statistics  
SCHUSS · Theory and Applications of Stochastic Differential Equations  
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  
SKINNER, HOLT, and SMITH · Analysis of Complex Surveys  
STOYAN, KENDALL, and MECKE · Stochastic Geometry and Its Applications, *Second Edition*  
STOYAN and STOYAN · Fractals, Random Shapes and Point Fields: Methods of Geometrical Statistics  
THOMPSON · Empirical Model Building  
THOMPSON · Sampling  
TIERNEY · LISP-STAT: An Object-Oriented Environment for Statistical Computing and Dynamic Graphics  
TIJMS · Stochastic Modeling and Analysis: A Computational Approach  
TITTERINGTON, SMITH, and MAKOV · Statistical Analysis of Finite Mixture Distributions  
UPTON and FINGLETON · Spatial Data Analysis by Example, Volume 1: Point Pattern and Quantitative Data  
UPTON and FINGLETON · Spatial Data Analysis by Example, Volume II: Categorical and Directional Data  
VAN RIJCKEVORSEL and DE LEEUW · Component and Correspondence Analysis  
WEISBERG · Applied Linear Regression, *Second Edition*  
WESTFALL and YOUNG · Resampling-Based Multiple Testing: Examples and Methods for  $p$ -Value Adjustment  
WHITTLE · Optimization Over Time: Dynamic Programming and Stochastic Control, Volume I and Volume II  
WHITTLE · Systems in Stochastic Equilibrium  
WONNACOTT and WONNACOTT · Econometrics, *Second Edition*  
WONNACOTT and WONNACOTT · Introductory Statistics, *Fifth Edition*  
WONNACOTT and WONNACOTT · Introductory Statistics for Business and Economics, *Fourth Edition*  
WOODING · Planning Pharmaceutical Clinical Trials: Basic Statistical Principles  
WOOLSON · Statistical Methods for the Analysis of Biomedical Data  
\*ZELLNER · An Introduction to Bayesian Inference in Econometrics
- Tracts on Probability and Statistics*  
BILLINGSLEY · Convergence of Probability Measures  
KELLY · Reversibility and Stochastic Networks  
TOUTENBURG · Prior Information in Linear Models

