

An Introduction to MULTIVARIATE STATISTICAL ANALYSIS 2nd Edition



T.W. Anderson

A Volume in the Wiley Series in Probability and Mathematical Statistics:
Vic Barnett, Ralph A. Bradley, J. Stuart Hunter, David G. Kendall,
Rupert G. Miller, Jr., Stephen M. Stigler, Geoffrey S. Watson
Advisory Editors

An Introduction to Multivariate Statistical Analysis

Second Edition

T. W. ANDERSON

Professor of Statistics and Economics
Stanford University

JOHN WILEY & SONS

New York Chichester Brisbane Toronto Singapore

Copyright © 1984 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.

Library of Congress Cataloging in Publication Data:

Anderson, T. W. (Theodore Wilbur), 1918–

An introduction to multivariate statistical analysis.

Bibliography: p.

Includes index.

1. Multivariate analysis. I. Title.

QA278.A516 1984 519.5'35 84-7334

ISBN 0-471-88987-3

Printed in the United States of America

10 9 8 7 6 5 4 3 2

Preface to the Second Edition

Twenty-six years have passed since the first edition of this book was published. During that time great advances have been made in multivariate statistical analysis—particularly in the areas treated in that volume. This new edition purports to bring the original edition up to date by substantial revision, rewriting, and additions. The basic approach has been maintained, namely, a mathematically rigorous development of statistical methods for observations consisting of several measurements or characteristics of each subject and a study of their properties. The general outline of topics has been retained.

The method of maximum likelihood has been augmented by other considerations. In point estimation of the mean vector and covariance matrix alternatives to the maximum likelihood estimators that are better with respect to certain loss functions, such as Stein and Bayes estimators, have been introduced. In testing hypotheses likelihood ratio tests have been supplemented by other invariant procedures. New results on distributions and asymptotic distributions are given; some significance points are tabulated. Properties of these procedures, such as power functions, admissibility, unbiasedness, and monotonicity of power functions, are studied. Simultaneous confidence intervals for means and covariances are developed. A chapter on factor analysis replaces the chapter sketching miscellaneous results in the first edition. Some new topics, including simultaneous equations models and linear functional relationships, are introduced. Additional problems present further results.

It is impossible to cover all relevant material in this book; what seems most important has been included. For a comprehensive listing of papers until 1966 and books until 1970 the reader is referred to *A Bibliography of Multivariate Statistical Analysis* by Anderson, Das Gupta, and Styan (1972). Further references can be found in *Multivariate Analysis: A Selected and Abstracted Bibliography, 1957–1972* by Subrahmaniam and Subrahmaniam (1973).

I am in debt to many students, colleagues, and friends for their suggestions and assistance; they include Yasuo Amemiya, James Berger, Byoung-Seon Choi, Arthur Cohen, Margery Cruise, Somesh Das Gupta, Kai-Tai Fang, Gene Golub, Aaron Han, Takeshi Hayakawa, Jogi Henna, Huang Hsu, Fred Huffer, Mituaki Huzii, Jack Kiefer, Mark Knowles, Sue Leurgans, Alex McMillan, Masashi No, Ingram Olkin, Kartik Patel, Michael Perlman, Allen Sampson, Ashis Sen Gupta, Andrew Siegel, Charles Stein, Patrick Strout, Akimichi Takemura, Joe Verducci, Marlos Viana, and Y. Yajima. I was helped in preparing the manuscript by Dorothy Anderson, Alice Lundin, Amy Schwartz, and Pat Struse. Special thanks go to Johanne Thiffault and George P. H. Styan for their precise attention. Support was contributed by the Army Research Office, the National Science Foundation, the Office of Naval Research, and IBM Systems Research Institute.

Seven tables of significance points are given in Appendix B to facilitate carrying out test procedures. Tables 1, 5, and 7 are Tables 47, 50, and 53, respectively, of *Biometrika Tables for Statisticians*, Vol. 2, by E. S. Pearson and H. O. Hartley; permission of the Biometrika Trustees is hereby acknowledged. Table 2 is made up from three tables prepared by A. W. Davis and published in *Biometrika* (1970a), *Annals of the Institute of Statistical Mathematics* (1970b), and *Communications in Statistics, B. Simulation and Computation* (1980). Tables 3 and 4 are Tables 6.3 and 6.4, respectively, of *Concise Statistical Tables*, edited by Ziro Yamauti (1977) and published by the Japanese Standards Association; this book is a concise version of *Statistical Tables and Formulas with Computer Applications*, JSA-1972. Table 6 is Table 3 of *The Distribution of the Sphericity Test Criterion*, ARL 72-0154, by B. N. Nagarsenker and K. C. S. Pillai, Aerospace Research Laboratories (1972). The author is indebted to the authors and publishers listed above for permission to reproduce these tables.

T. W. ANDERSON

Stanford, California
June 1984

Preface to the First Edition

This book has been designed primarily as a text for a two-semester course in multivariate statistics. It is hoped that the book will also serve as an introduction to many topics in this area to statisticians who are not students and will be used as a reference by other statisticians.

For several years the book in the form of dittoed notes has been used in a two-semester sequence of graduate courses at Columbia University; the first six chapters constituted the text for the first semester, emphasizing correlation theory. It is assumed that the reader is familiar with the usual theory of univariate statistics, particularly methods based on the univariate normal distribution. A knowledge of matrix algebra is also a prerequisite; however, an appendix on this topic has been included.

It is hoped that the more basic and important topics are treated here, though to some extent the coverage is a matter of taste. Some of the more recent and advanced developments are only briefly touched on in the last chapter.

The method of maximum likelihood is used to a large extent. This leads to reasonable procedures; in some cases it can be proved that they are optimal. In many situations, however, the theory of desirable or optimum procedures is lacking.

Over the years this manuscript has been developed, a number of students and colleagues have been of considerable assistance. Allan Birnbaum, Harold Hotelling, Jacob Horowitz, Howard Levene, Ingram Olkin, Gobind Seth, Charles Stein, and Henry Teicher are to be mentioned particularly. Acknowledgments are also due to other members of the Graduate Mathematical

Statistics Society at Columbia University for aid in the preparation of the manuscript in dittoed form. The preparation of this manuscript was supported in part by the Office of Naval Research.

T. W. ANDERSON

*Center for Advanced Study
in the Behavioral Sciences
Stanford, California
December 1957*

Contents

CHAPTER 1	
Introduction	1
1.1. Multivariate Statistical Analysis	1
1.2. The Multivariate Normal Distribution	3
CHAPTER 2	
The Multivariate Normal Distribution	6
2.1. Introduction	6
2.2. Notions of Multivariate Distributions	7
2.3. The Multivariate Normal Distribution	14
2.4. The Distribution of Linear Combinations of Normally Distributed Variates; Independence of Variates; Marginal Distributions	24
2.5. Conditional Distributions and Multiple Correlation Coefficient	35
2.6. The Characteristic Function; Moments Problems	43 50
CHAPTER 3	
Estimation of the Mean Vector and the Covariance Matrix	59
3.1. Introduction	59
3.2. The Maximum Likelihood Estimators of the Mean Vector and the Covariance Matrix	60
3.3. The Distribution of the Sample Mean Vector; Inference Concerning the Mean When the Covariance Matrix Is Known	68

3.4. Theoretical Properties of Estimators of the Mean Vector	77
3.5. Improved Estimation of the Mean Problems	86 96
CHAPTER 4	
The Distributions and Uses of Sample Correlation Coefficients	102
4.1. Introduction	102
4.2. Correlation Coefficient of a Bivariate Sample	103
4.3. Partial Correlation Coefficients; Conditional Distributions	125
4.4. The Multiple Correlation Coefficient Problems	134 149
CHAPTER 5	
The Generalized T^2 -Statistic	156
5.1. Introduction	156
5.2. Derivation of the Generalized T^2 -Statistic and Its Distribution	157
5.3. Uses of the T^2 -Statistic	164
5.4. The Distribution of T^2 Under Alternative Hypotheses; The Power Function	173
5.5. The Two-Sample Problem with Unequal Covariance Matrices	175
5.6. Some Optimal Properties of the T^2 -Test Problems	181 190
CHAPTER 6	
Classification of Observations	195
6.1. The Problem of Classification	195
6.2. Standards of Good Classification	196
6.3. Procedures of Classification into One of Two Populations with Known Probability Distributions	199
6.4. Classification into One of Two Known Multivariate Normal Populations	204
6.5. Classification into One of Two Multivariate Normal Populations When the Parameters Are Estimated	208
6.6. Probabilities of Misclassification	217
6.7. Classification into One of Several Populations	224

6.8. Classification into One of Several Multivariate Normal Populations	228
6.9. An Example of Classification into One of Several Multivariate Normal Populations	231
6.10. Classification into One of Two Known Multivariate Normal Populations with Unequal Covariance Matrices	234
Problems	241

CHAPTER 7

The Distribution of the Sample Covariance Matrix and the Sample Generalized Variance	244
7.1. Introduction	244
7.2. The Wishart Distribution	245
7.3. Some Properties of the Wishart Distribution	252
7.4. Cochran's Theorem	257
7.5. The Generalized Variance	259
7.6. Distribution of the Set of Correlation Coefficients When the Population Covariance Matrix Is Diagonal	266
7.7. The Inverted Wishart Distribution and Bayes Estimation of the Covariance Matrix	268
7.8. Improved Estimation of the Covariance Matrix	273
Problems	279

CHAPTER 8

Testing the General Linear Hypothesis; Multivariate Analysis of Variance	285
8.1. Introduction	285
8.2. Estimators of Parameters in Multivariate Linear Regression	287
8.3. Likelihood Ratio Criteria for Testing Linear Hypotheses About Regression Coefficients	292
8.4. The Distribution of the Likelihood Ratio Criterion When the Hypothesis Is True	298
8.5. An Asymptotic Expansion of the Distribution of the Likelihood Ratio Criterion	311
8.6. Other Criteria for Testing the Linear Hypothesis	321
8.7. Tests of Hypotheses About Matrices of Regression Coefficients and Confidence Regions	333

8.8. Testing Equality of Means of Several Normal Distributions with Common Covariance Matrix	338
8.9. Multivariate Analysis of Variance	342
8.10. Some Optimal Properties of Tests Problems	349
CHAPTER 9	
Testing Independence of Sets of Variates	376
9.1. Introduction	376
9.2. The Likelihood Ratio Criterion for Testing Independence of Sets of Variates	376
9.3. The Distribution of the Likelihood Ratio Criterion When the Null Hypothesis Is True	381
9.4. An Asymptotic Expansion of the Distribution of the Likelihood Ratio Criterion	385
9.5. Other Criteria	387
9.6. Step-down Procedures	389
9.7. An Example	392
9.8. The Case of Two Sets of Variates	394
9.9. Admissibility of the Likelihood Ratio Test	397
9.10. Monotonicity of Power Functions of Tests of Independence of Sets Problems	399
CHAPTER 10	
Testing Hypotheses of Equality of Covariance Matrices and Equality of Mean Vectors and Covariance Matrices	404
10.1. Introduction	404
10.2. Criteria for Testing Equality of Several Covariance Matrices	405
10.3. Criteria for Testing That Several Normal Distributions Are Identical	408
10.4. Distributions of the Criteria	410
10.5. Asymptotic Expansions of the Distributions of the Criteria	419
10.6. The Case of Two Populations	422
10.7. Testing the Hypothesis That a Covariance Matrix Is Proportional to a Given Matrix; The Sphericity Test	427

10.8. Testing the Hypothesis That a Covariance Matrix Is Equal to a Given Matrix	434
10.9. Testing the Hypothesis That a Mean Vector and a Covariance Matrix Are Equal to a Given Vector and Matrix	440
10.10. Admissibility of Tests	443
Problems	446

CHAPTER 11

Principal Components	451
11.1. Introduction	451
11.2. Definition of Principal Components in the Population	452
11.3. Maximum Likelihood Estimators of the Principal Components and Their Variances	460
11.4. Computation of the Maximum Likelihood Estimates of the Principal Components	462
11.5. An Example	465
11.6. Statistical Inference	468
11.7. Testing Hypotheses about the Characteristic Roots of a Covariance Matrix	473
Problems	477

CHAPTER 12

Canonical Correlations and Canonical Variables	480
12.1. Introduction	480
12.2. Canonical Correlations and Variates in the Population	481
12.3. Estimation of Canonical Correlations and Variates	492
12.4. Statistical Inference	497
12.5. An Example	500
12.6. Linearly Related Expected Values	502
12.7. Simultaneous Equations Models	509
Problems	519

CHAPTER 13

The Distributions of Characteristic Roots and Vectors	521
13.1. Introduction	521
13.2. The Case of Two Wishart Matrices	522

13.3.	The Case of One Nonsingular Wishart Matrix	532
13.4.	Canonical Correlations	538
13.5.	Asymptotic Distributions in the Case of One Wishart Matrix	540
13.6.	Asymptotic Distributions in the Case of Two Wishart Matrices Problems	544 548
CHAPTER 14		
	Factor Analysis	550
14.1.	Introduction	550
14.2.	The Model	551
14.3.	Maximum Likelihood Estimators for Random Orthogonal Factors	557
14.4.	Estimation for Fixed Factors	569
14.5.	Factor Interpretation and Transformation	570
14.6.	Estimation for Identification by Specified Zeros	574
14.7.	Estimation of Factor Scores Problems	575 576
APPENDIX A		
	Matrix Theory	579
A.1.	Definition of a Matrix and Operations on Matrices	579
A.2.	Characteristic Roots and Vectors	587
A.3.	Partitioned Vectors and Matrices	591
A.4.	Some Miscellaneous Results	596
A.5.	Gram-Schmidt Orthogonalization and the Solution of Linear Equations	605
APPENDIX B		
	Tables	609
1.	Wilks' Likelihood Criterion: Factors $C(p, m, M)$ to Adjust to χ^2_{pm} where $M = n - p + 1$	609
2.	Tables of Significance Points for the Lawley-Hotelling Trace Test	616
3.	Tables of Significance Points for the Bartlett-Nanda-Pillai Trace Test	630
4.	Tables of Significance Points for the Roy Maximum Root Test	634

CONTENTS

xvii

5. Tables of Significance Points for the Modified Likelihood Ratio Test of Equality of Covariance Matrices Based on Equal Sample Sizes	638
6. Correction Factors for Significance Points for the Sphericity Test	639
7. Significance Points for the Modified Likelihood Ratio Test $\Sigma = \Sigma_0$	641
References	643
Index	667

CHAPTER 1

Introduction

1.1. MULTIVARIATE STATISTICAL ANALYSIS

Multivariate statistical analysis is concerned with data that consist of sets of measurements on a number of individuals or objects. The sample data may be heights and weights of some individuals drawn randomly from a population of schoolchildren in a given city, or the statistical treatment may be made on a collection of measurements, such as lengths and widths of petals and lengths and widths of sepals of iris plants taken from two species, or one may study the scores on batteries of mental tests administered to a number of students.

The measurements made on a single individual can be assembled into a column vector. We think of the entire vector as an observation from a multivariate population or distribution. When the individual is drawn randomly, we consider the vector as a random vector with a distribution or probability law describing that population. The set of observations on all individuals in a sample constitute a sample of vectors, and the vectors set side by side make up the matrix of observations.[†] The data to be analyzed then are thought of as displayed in a matrix or in several matrices.

We shall see that it is helpful in visualizing the data and understanding the methods to think of each observation vector as constituting a point in a Euclidean space, each coordinate corresponding to a measurement or variable. Indeed, an early step in the statistical analysis is plotting the data; since most statisticians are limited to two-dimensional plots, two coordinates of the observation are plotted in turn.

[†] When data are listed on paper by individual, it is natural to print the measurements on one individual as a row of the table; then one individual corresponds to a *row* vector. Since we prefer to operate algebraically with column vectors, we have chosen to treat observations in terms of *column* vectors. (In practice, the basic data set may well be on cards, tapes, or disks.)

Characteristics of a univariate distribution of essential interest are the mean as a measure of location and the standard deviation as a measure of variability; similarly the mean and standard deviation of a univariate sample are important summary measures. In multivariate analysis, the means and variances of the separate measurements—for distributions and for samples—have corresponding relevance. An essential aspect, however, of multivariate analysis is the dependence between the different variables. The dependence between two variables may involve the covariance between them, that is, the average products of their deviations from their respective means. The covariance standardized by the corresponding standard deviations is the correlation coefficient; it serves as a measure of degree of dependence. A set of summary statistics is the mean vector (consisting of the univariate means) and the covariance matrix (consisting of the univariate variances and bivariate covariances). An alternative set of summary statistics with the same information is the mean vector, the set of standard deviations, and the correlation matrix. Similar parameter quantities describe location, variability, and dependence in the population or for a probability distribution. The multivariate *normal* distribution is completely determined by its mean vector and covariance matrix, and the sample mean vector and covariance matrix constitute a sufficient set of statistics.

The measurement and analysis of dependence between variables, between sets of variables, and between variables and sets of variables, are fundamental to multivariate analysis. The multiple correlation coefficient is an extension of the notion of correlation to the relationship of one variable to a set of variables. The partial correlation coefficient is a measure of dependence between two variables when the effects of other correlated variables have been removed. The various correlation coefficients computed from samples are used to estimate corresponding correlation coefficients of distributions. In this book tests of hypotheses of independence are developed. The properties of the estimators and test procedures are studied for sampling from the multivariate normal distribution.

A number of statistical problems arising in multivariate populations are straightforward analogs of problems arising in univariate populations; the suitable methods for handling these problems are similarly related. For example, in the univariate case we may wish to test the hypothesis that the mean of a variable is zero; in the multivariate case we may wish to test the hypothesis that the vector of the means of several variables is the zero vector. The analog of the Student *t*-test for the first hypothesis is the generalized T^2 -test. The

analysis of variance of a single variable is adapted to vector observations; in regression analysis, the dependent quantity may be a vector variable. A comparison of variances is generalized into a comparison of covariance matrices.

The test procedures of univariate statistics are generalized to the multivariate case in such ways that the dependence between variables is taken into account. These methods may not depend on the coordinate system; in other terms, the procedures are invariant with respect to linear transformations that leave the null hypothesis invariant. In some problems there may be families of tests that are invariant; then choices must be made. Optimal properties of the tests are considered.

For some other purposes, however, it may be important to select a coordinate system so that the variates have desired statistical properties. One might say that they involve characterizations of inherent properties of normal distributions and of samples. These are closely related to the algebraic problems of canonical forms of matrices. An example is finding the normalized linear combination of variables with maximum or minimum variance (finding principal components); this amounts to finding a rotation of axes that carries the covariance matrix to diagonal form. Another example is characterizing the dependence between two sets of variates (finding canonical correlations). These problems involve the characteristic roots and vectors of various matrices. The statistical properties of the corresponding sample quantities are treated.

Some statistical problems arise in models in which means and covariances are restricted. Factor analysis may be based on a model with a (population) covariance matrix that is the sum of a positive definite diagonal matrix and a positive semidefinite matrix of low rank; linear structural relationships may have a similar formulation. The simultaneous equations system of econometrics is another example of a special model.

1.2. THE MULTIVARIATE NORMAL DISTRIBUTION

The statistical methods treated in this book can be developed and evaluated in the context of the multivariate normal distribution, though many of the procedures are useful and effective when the distribution sampled is not normal. A major reason for basing statistical analysis on the normal distribution is that this probabilistic model approximates well the distribution of continuous measurements in many sampled populations. In fact, most of the