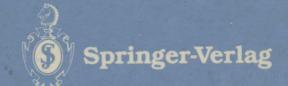
Ira H. Bernstein

Applied Multivariate Analysis

With Calvin P. Garbin and Gary K. Teng



Ira H. Bernstein

Applied Multivariate Analysis

With Calvin P. Garbin and Gary K. Teng

With 37 Illustrations



Springer-Verlag New York Berlin Heidelberg London Paris Tokyo Ira H. Bernstein Department of Psychology University of Texas at Arlington Arlington, Texas 76019 USA

Calvin P. Garbin Department of Psychology University of Nebraska at Lincoln Lincoln, Nebraska 68588 USA

Gary K. Teng Technical Evaluation and Management Systems, Inc. (TEAMS®) Dallas, Texas 75240 LISA

Library of Congress Cataloging-in-Publication Data Bernstein, Ira H.

Applied multivariate analysis.

Bibliography: p.

Includes index.

1. Multivariate analysis. I. Garbin, Calvin P.

II. Teng, Garv K. III. Title.

QA278.B457 1987 519.5'35 87-15174

© 1988 by Springer-Verlag New York Inc.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer-Verlag, 175 Fifth Avenue, New York, New York 10010, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden. The use of general descriptive names, trade names, trademarks, etc. in this publication, even if the former are not especially identified, is not to be taken as a sign that such names, as understood by the Trade Marks and Merchandise Marks Act, may accordingly be used freely by anyone.

Typeset by Asco Trade Typesetting Ltd., Hong Kong. Printed and bound by R.R. Donnelley & Sons, Harrisonburg, Virginia. Printed in the United States of America.

987654321

ISBN 0-387-96542-4 Springer-Verlag New York Berlin Heidelberg ISBN 3-540-96542-4 Springer-Verlag Berlin Heidelberg New York

To Linda, Cari, Dina and the memory of Jum C. Nunnally

Preface

Like most academic authors, my views are a joint product of my teaching and my research. Needless to say, my views reflect the biases that I have acquired. One way to articulate the rationale (and limitations) of my biases is through the preface of a truly great text of a previous era, Cooley and Lohnes (1971, p. v). They draw a distinction between mathematical statisticians whose intellect gave birth to the field of multivariate analysis, such as Hotelling, Bartlett, and Wilks, and those who chose to "concentrate much of their attention on methods of analyzing data in the sciences and of interpreting the results of statistical analysis (and) ... who are more interested in the sciences than in mathematics, among other characteristics."

I find the distinction between individuals who are temperamentally "mathematicians" (whom philosophy students might call "Platonists") and "scientists" ("Aristotelians") useful as long as it is not pushed to the point where one assumes "mathematicians" completely disdain data and "scientists" are never interested in contributing to the mathematical foundations of their discipline. I certainly feel more comfortable attempting to contribute in the "scientist" rather than the "mathematician" role.

As a consequence, this book is primarily written for individuals concerned with data analysis. However, as noted in Chapter 1, true expertise demands familiarity with both traditions.

One consequence of my bias is that even though I have a great love for data, I have long since learned not to worship a particular data set, i.e., I believe in sampling error and replication. I especially believe in Henry Kaiser's (1970) aphorism "It don't make no nevermind" about highly elaborate weighting schemes that more often than not prove less useful than simpler ones such as weighting variables equally. (Kaiser, by the way, shows the limitations of the "rule of thumb" distinction between Mathematician/Platonists and Scientist/Aristotelians; bright people do both as his numerous other contributions such as the varimax rotation and alpha factoring attest.) If you are unsure as to what I mean by "magical equations," pick up almost any scholarly journal containing path analyses and confirmatory factor analyses. Similarly, it should not be too difficult to locate a book in which the author attempts

to use the force of copyright law to protect the weighting scheme used in an equation.

The major reason that I refer to Cooley and Lohnes (1971) as belonging to a prior generation is that it was written before the ascendancy of the major computer packages. I assume you will be using these packages. At one point in writing this book, I thought of offering a "translation" table to explain the printouts that are provided by the packages. The reason that I did not devote more space to the topic is the rate at which revisions of the packages make such information obsolete.

The one feature that I hope sets this book apart from the many other excellent books in the field is my use of the computer as a device to teach you about data structures. It is all too easy to see only the data analytic features of computer packages, as wonderful as they are. (Consider that in almost no time at all I can enter the coding to do all variants on a particular analysis that would have been prohibitive in terms of time when I was a student.) Much less apparent to the student is the way that a computer can *generate* data to conform to a model or deviate in specified ways. Scholars have long used Monte Carlo and related simulation approaches; I feel that it is sufficiently easy for students at this level, or even at a more introductory level, to perform simulations so that they should be part of all students' training.

The material in this book is designed to be covered in a standard one semester graduate course. It is assumed that a student has had a conventional first semester graduate class in statistics. I have, however, devoted a chapter to material that is largely a summary of this core material (Chapter 2). Of course, I have also tried to include useful information that would make the book valuable beyond the formal confines of the course (*lagniappe*, as one would say in the language of those to my Southeast).

Acknowledgments

I would like to thank Springer-Verlag for their faith in this project. Robert Widner's comments on various points of the text were most highly appreciated. Various ideas were outgrowths of conversations I had with various departmental colleagues, especially Jim Nairne, Jim Erickson, Jim Baerwaldt, Bill Ickes, Paul Paulus, Duane Martin, and the ophthalmologists at the Glaucoma Associates of Texas. I also am grateful to Melvin Pierce, Tom Kennedy, Ed Homko, Michael Griffin, and others at the University of Texas Academic Computer center, Rory Gresham of Seabrook Computers, and Sara Kapka of IBM for their technical assistance. The oblique, but nonetheless vital contributions of Emma Peel, Ferdinand LaMenthe, John B. Gillespie, John Coltrane, Charles C. Parker, Jr., Larson E. Whipsnade, Jonathan Steed, John Sheridan, and James Cullum and associates are also noted. Finally, the wisdom of a comment by my long-term friend and colleague, Professor Stanley

Coren of the University of British Columbia that "you never finish writing a book, they just take it away from you" is hereby duly noted.

Arlington, Texas May 1987 IRA H. BERNSTEIN

Contents

| Preface | |
|-------------------------------------------------|--|
| CHAPTER 1 INTRODUCTION AND PREVIEW | |
| Chapter Overview | |
| Multivariate Analysis: A Broad Definition | |
| Multivariate Analysis: A Narrow Definition | |
| Some Important Themes | |
| Obtaining Meaningful Relations | |
| Selecting Cutoffs | |
| Questions of Statistical Inference | |
| Outliers | |
| The Importance of Theory | |
| Problems Peculiar to the Analysis of Scales | |
| The Role of Computers in Multivariate Analysis | |
| Multivariate Analysis and the Personal Computer | |
| Choosing a Computer Package | |
| Problems in the Use of Computer Packages | |
| The Importance of Matrix Procedures | |
| | |
| CHAPTER 2 SOME BASIC STATISTICAL CONCEPTS | |
| Chapter Overview | |
| Univariate Data Analysis | |
| Frequency Distributions | |
| Normal Distributions | |
| Standard Normal Distributions | |
| Parameters and Statistics | |
| Locational Parameters and Statistics | |
| Measures of Variability | |
| A Note on Estimation | |
| Binary Data and the Binomial Distribution | |
| Data Transformation | |
| Bivariate Data Analysis | |
| Characteristics of Bivariate Relationships | |
| Bivariate Normality | |

| Measures of Bivariate Relation | 38 |
|---------------------------------------------------------------|-----|
| Range Restriction | 39 |
| Pearson Correlation Formulas in Special Cases | 40 |
| Non-Pearson Estimates of Pearson Correlations | 40 |
| The Eta-Square Measure | 41 |
| Phi Coefficients with Unequal Probabilities | 42 |
| Sampling Error of a Correlation | 42 |
| The Z' Transformation | 42 |
| Linear Regression | 43 |
| The Geometry of Regression | 43 |
| Raw-Score Formulas for the Slope | 45 |
| Raw-Score Formulas for the Intercept | 45 |
| Residuals | 45 |
| The Standard Error of Estimate | 46 |
| Why the Term "Regression"? | 48 |
| A Summary of Some Basic Relations | 49 |
| Statistical Control: A First Look at Multivariate Relations | 49 |
| Partial and Part Correlation | 49 |
| Statistical versus Experimental Control | 53 |
| Multiple Partialling | 53 |
| Within-Group, Between-Group, and Total Correlations | 54 |
| | |
| CHAPTER 3 SOME MATRIX CONCEPTS | 57 |
| CHAFTER 5 SOME MATRIX CONCEPTS | 37 |
| Chapter Overview | 57 |
| Basic Definitions | 59 |
| Square Matrices | 61 |
| Transposition | 64 |
| Matrix Equality | 64 |
| Basic Matrix Operations. | 65 |
| Matrix Addition and Subtraction | 65 |
| Matrix Multiplication | 66 |
| Correlation Matrices and Matrix Multiplication | 68 |
| Partitioned Matrices and Their Multiplication | 69 |
| Some Rules and Theorems Involved in Matrix Algebra | 69 |
| Products of Symmetric Matrices | 70 |
| More about Vector Products | 71 |
| Exponentiation | 71 |
| Determinants | 72 |
| Matrix Singularity and Linear Dependency | 73 |
| Matrix Rank | 74 |
| Matrix "Division" | 75 |
| The Inverse of a 2 × 2 Matrix | 75 |
| Inverses of Higher-Order Matrices | 75 |
| Recalculation of an Inverse Following Deletion of Variable(s) | 76 |
| An Application of Matrix Algebra | 77 |
| More about Linear Combinations | 79 |
| The Mean of a Linear Combination | 80 |
| The Variance of a Linear Combination | 80 |
| Covariances between Linear Combination | Q 1 |

| | Contents |
|-------------------|----------------------------------------------------------------|
| The Correlation | on between Two Different Linear Combinations |
| | between Linear Combinations and Matrix Notation |
| | |
| | Make No Nevermind" Principle |
| | Eigenvectors |
| | nanalysis |
| Eigenanalysis | of Gramian Matrices |
| CHAPTER 4 | MULTIPLE REGRESSION AND |
| | CORRELATION—PART 1. BASIC CONCEPTS |
| Chapter Overvie | |
| | derlying Multiple Regression |
| _ | ate Normal Distribution |
| | cometry of Multiple Regression |
| | egression Analysis |
| | Predictors |
| | |
| | nple |
| | ppressor Variables |
| • | d Formulas |
| | rmulas |
| | ons for R^2 |
| | he Relative Importance of the Two Predictors |
| | le Correlation |
| | re Than Two Predictors |
| | Multicollinearity |
| Another Way | to Obtain R ² |
| Residuals | |
| Inferential Tests | |
| Testing $R \dots$ | |
| Testing Beta V | Veights |
| | niqueness of Predictors |
| | native Equations |
| | on |
| | Correlation from a priori Weights |
| | fference between R^2 and r^2 Derived from a priori Weights |
| | nclusion of Predictors |
| | sion of Predictors |
| | Handle Multicollinearity |
| | |
| Comparing Al | ternative Equations |
| | fect Prediction. |
| Example 2—Imp | perfect Prediction plus a Look at Residuals |
| | l Personality Assessment Data |
| Alternative Appr | roaches to Data Aggregation |
| CHAPTER 5 | MULTIPLE REGRESSION AND |
| | CORRELATION—PART 2. |
| | ADVANCED APPLICATIONS |
| | |
| • | W |
| Nonguantitative | Variables |

xiv Contents

| Dummy Coding | 123 |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|-----------------------------------------------------------------------------------------------|
| Effect Coding | 124 |
| Orthogonal Coding | 120 |
| The Simple Analysis of Variance (ANOVA) | 12 |
| Fixed Effects versus Random Effects | 127 |
| The Simple ANOVA as a Formal Model | 128 |
| Results of Regression ANOVAs | 131 |
| Multiple Comparisons | 133 |
| Orthogonal versus Nonorthogonal Contrasts | 133 |
| Planned versus Unplanned Comparisons | 134 |
| Individual Alpha Levels versus Groupwise Alpha Levels | 135 |
| Evaluation of Quantitative Relations | 135 |
| Method I | 136 |
| Method II. | 138 |
| The Two-Way ANOVA | 138 |
| Equal-N Analysis | 142 |
| Unequal-N Analysis | 143 |
| Fitting Parallel Lines | 145 |
| Simple Effect Models | 148 |
| The Analysis of Covariance (ANCOVA) | 150 |
| Effects of the ANCOVA on the Treatment Sum of Squares | 15 |
| | 15 |
| Using Dummy Codes to Plot Group Means | |
| Repeated Measures, Blocked and Matched Designs | 15. |
| Higher-Order Designs | 150 |
| CHAPTER 6 EXPLORATORY FACTOR ANALYSIS | 15 |
| Chapter Overview | 15 |
| Chapter Overview | 15° |
| Chapter Overview | 15′ 16 16 |
| Chapter Overview | 15′ 16 16 |
| Chapter Overview | 15 ⁷ 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. | 15' 16' 16' 16' 16' |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix. | 15' 16' 16' 16' 16' |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix | 15° 16° 16° 16° 16° 16° 16° |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix Error Scores and Error Loadings. | 15° 16° 16° 16° 16° 16° 16° |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. | 15 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. | 15 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation The First Form of Factor Indeterminacy. An Important Special Case | 15 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix Error Scores and Error Loadings. The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. | 15 ⁷ 16 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix Error Scores and Error Loadings. The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization | 15 ⁷ 16 16 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization Reduction in the Number of Variables | 15 ⁷ 16 16 16 16 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix Factor Scores and the Factor Score Matrix Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization Reduction in the Number of Variables Dimensional Analysis. | 15° 166 166 166 166 166 166 166 166 166 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization Reduction in the Number of Variables Dimensional Analysis Determination of Factor Scores An Overview of the Exploratory Factoring Process | 15° 166 166 166 166 166 166 166 166 166 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization Reduction in the Number of Variables Dimensional Analysis Determination of Factor Scores An Overview of the Exploratory Factoring Process Principal Components | 15° 16° 16° 16° 16° 16° 16° 16° 16° 16° 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix. Error Scores and Error Loadings. The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization. Reduction in the Number of Variables. Dimensional Analysis. Determination of Factor Scores. An Overview of the Exploratory Factoring Process. Principal Components. The Eigenvectors and Eigenvalues of a Gramian Matrix. | 15° 16° 16° 16° 16° 16° 16° 16° 16° 16° 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix. Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization. Reduction in the Number of Variables. Dimensional Analysis. Determination of Factor Scores. An Overview of the Exploratory Factoring Process. Principal Components. The Eigenvectors and Eigenvalues of a Gramian Matrix. A Note on the Orthogonality of PCs. | 15 16 16 16 16 16 16 16 16 16 16 16 16 16 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix. Error Scores and Error Loadings. The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization. Reduction in the Number of Variables. Dimensional Analysis. Determination of Factor Scores. An Overview of the Exploratory Factoring Process. Principal Components. The Eigenvectors and Eigenvalues of a Gramian Matrix. A Note on the Orthogonality of PCs. How an Eigenanalysis Is Performed. | 15 16 16 16 16 16 16 16 16 16 16 16 16 17 17 |
| Chapter Overview. The Basic Factor Analytic Model. The Factor Equation. The Raw-Score Matrix. Factor Scores and the Factor Score Matrix. Pattern Elements and the Pattern Matrix. Error Scores and Error Loadings The Covariance Equation. The First Form of Factor Indeterminacy. An Important Special Case. Common Uses of Factor Analysis. Orthogonalization. Reduction in the Number of Variables. Dimensional Analysis. Determination of Factor Scores. An Overview of the Exploratory Factoring Process. Principal Components. The Eigenvectors and Eigenvalues of a Gramian Matrix. A Note on the Orthogonality of PCs. | 15° 15° 16° 16° 16° 16° 16° 16° 16° 16° 16° 16 |

| | Contents |
|----------------------------------------------------------------------------------------------------------------------|----------|
| Factor Definition and Rotation | |
| Factor Definition | |
| Simple Structure | |
| PC versus Simple Structure | |
| Graphic Representation | |
| Analytic Orthogonal Rotation | |
| Oblique Rotations | |
| Reference Vectors | |
| Analytic Oblique Rotation | |
| The Common Factor Model | |
| An Example of the Common Factor Model | |
| A Second Form of Factor Indeterminacy | |
| Factor Scores | |
| "Exact" Procedures | |
| Estimation Procedures | |
| Approximation Procedures | |
| Addendum: Constructing Correlation Matrices with a Desired Fac | |
| Structure | |
| | |
| CHAPTER 7 CONFIRMATORY FACTOR ANALYSI | S |
| | |
| Chapter Overview | |
| Comparing Factor Structures | |
| Similarity of Individual Factors versus Similarity of the Overall | |
| Solution | |
| Case I—Comparing Alternate Solutions Derived from the Same Case II—Comparing Solutions Obtained from the Same Subjec | |
| Different Variables | |
| Case III—Comparing Solutions with the Same Variables but on | |
| Individuals; Matrix Information Available | |
| Case IV—Comparing Solutions with the Same Variables but Di | |
| Individuals; Matrix Information Unavailable | |
| Case V—Factor Matching | |
| Oblique Multiple Groups Tests of Weak Structure | |
| Basic Approach | |
| Evaluating the Substantive Model. | |
| Alternative Models | |
| Performing the Actual OMG Analysis | |
| Computational Steps | |
| OMG Common Factor Solutions | |
| A Numerical Example | |
| LISREL Tests of Weak Substantive Models | |
| Specification of Weak Models | |
| Estimation | |
| Identification | |
| Assessment of Fit | |
| Numerical Examples | |
| LISREL Tests of Strong Substantive Models | |
| Causal Models and Path Analysis. | |
| Causai Models and I am Analysis | |

| Example I: Or Example II: T Example III: O Decomposing Causal Models a A Formal Sep The Confirma The Confirma The Structura | Path Analytic Concepts | 228 231 233 234 236 237 237 238 239 |
|-----------------------------------------------------------------------------------------------------------------------------------------------------|----------------------------------------------------------------------------------------------|-------------------------------------------------------------|
| CHAPTER 8 | CLASSIFICATION METHODS—PART 1. FORMING DISCRIMINANT AXES | 246 |
| Chapter Overvie | ew | 246 |
| | nalysis with Two Groups and Two Predictors | 247 |
| | esentation | 248 |
| | Obtain Discriminant Weights | 250 |
| Another Way | to Obtain Discriminant Weights | 251 |
| | | 252 |
| | to Obtain Discriminant Weights | 256 |
| | of $\mathbf{W}^{-1}\mathbf{B}$ versus Eigenanalysis of $\mathbf{SP}_{w}^{-1}\mathbf{SP}_{b}$ | 257 |
| | elation | 257 |
| | Scores, Discriminant Means, and the Mahalanobis | |
| | sure | 258 |
| | ant Structure | 259 |
| | Discriminant Axis | 259 |
| | alysis with Two Predictors and Three Groups | 260 |
| | ultiple Discriminant Axes | 262 |
| | centration from the Eigenanalysis | 263 |
| | amples | 263 |
| Depiction of C | Group Differences in the Diffuse Example | 269 |
| | tion | 271 |
| | Veights | 271 |
| | icity | 272 |
| | alysis—The General Case | 274 |
| | as Discriminant Variables | 274 |
| | esting | 274 |
| | Factor Analysis | 275 |
| | riminant Analysis | 275 |
| CHAPTER 9 | CLASSIFICATION METHODS—PART 2. | |
| | METHODS OF ASSIGNMENT | 276 |
| Chanter Overvie | w | 276 |
| | ance Gaussian Model | 279 |
| | s and Cutoff Rules | 280 |
| Other Dresent | | 200 |

| | Contents |
|----------------------------------------------------|-----------|
| Why a Cutoff Rule? | |
| | |
| Bayesian Considerations | |
| Why Bayesian Considerations Are Important | |
| Varying Base Rates | |
| Bayes' Theorem in Binary Classification | |
| Receiver (Relative) Operating Characteristic (RC | |
| Describing Accuracy of Classification (Sensitivity | |
| The Unequal Variance Gaussian Model | |
| Other Signal Detection Models | ********* |
| Strategies for Individual Classification | ********* |
| Describing Performance | |
| Bayesian Considerations with Multiple Groups. | |
| Alternative Strategies—An Overview | |
| A Numerical Example | |
| Calibration Study | |
| The Context of Classification | |
| Classification Based on Salient Variables | |
| Homoscedasticity Assumed | |
| Normality Assumed | |
| Dealing with Heteroscedasticity | |
| Discriminant Functions and Classification | |
| | |
| Homoscedasticity Assumed | |
| Normality Assumed | |
| Dealing with Heteroscedasticity | |
| Using Multiple Discriminant Scores | |
| Normality and Multiple Discriminant Functions | 8 |
| Classification Based on Distance Measures | |
| Simple Distance | |
| Fisher's Classification Functions | |
| Mahalanobis Distances | |
| A Summary of Strategic Considerations in Classific | cation |
| CHAPTER 10 CLASSIFICATION METH | ODC DARTS |
| | |
| INFERENTIAL CONSIDE | |
| MANOVA | ***** |
| Chapter Overview | |
| Chapter Overview | |
| The Two-Group MANOVA and Hotelling's T^2 | |
| The Formal MANOVA Model | |
| Hotelling's T^2 | |
| Post Hoc Comparisons | |
| Application of Hotelling's T^2 to a Single Group | |
| Testing for Equality of Means | |
| Single Group MANOVA versus Repeated Meas | |
| Tests of Vector Means with Multiple Groups | |
| The Fundamental Problem | |
| Testing a Concentrated Structure | |
| Testing a Diffuse Structure | |
| Testing for Homogeneity of Covariance | |
| 0 | |

xviii Contents

| The Simple MANOVA with Multiple Groups | 332 |
|-------------------------------------------------------------------------|------|
| SP and Variance-Covariance Matrices | 333 |
| Applying Box's M Test | 334 |
| Specifying the Nature of the Group Differences | 335 |
| The Multivariate MANOVA | 336 |
| Terminology and Basic Logic | 336 |
| A Numerical Example | 338 |
| The A Effect | 339 |
| The B Effect | 340 |
| The AB Interaction | 340 |
| The MANCOVA | 341 |
| Numerical Example | 342 |
| | |
| CHAPTER 11 PROFILE AND CANONICAL ANALYSIS | 345 |
| Chapter Overview | 345 |
| Profile Similarity | 346 |
| Similarity of Scalars | 346 |
| Similarity of Vectors | 347 |
| Measuring Profile Elevation | 349 |
| Profile Similarity Based on Elevation Alone | 351 |
| Numerical Example | 351 |
| Profile Similarity Based on Shape Information | 353 |
| A Numerical Example | 354 |
| Correcting for Elevation | 355 |
| Correlations, Covariances, and Cross Products as Alternative Indices of | 333 |
| Similarity | 355 |
| | 357 |
| Simple and Hierarchical Clustering. | 360 |
| Numerical Example | |
| Canonical Analysis | 363 |
| Goals of Canonical Analysis | 363 |
| Basic Logic | 363 |
| Redundancy Analysis | 365 |
| Basic Matrix Operations | 366 |
| Statistical Inference | 369 |
| Numerical Example | 370 |
| Alternatives to Canonical Analysis | 375 |
| CHAPTER 12 ANALYSIS OF SCALES | 376 |
| Chapter Overview | 376 |
| Properties of Individual Items. | 380 |
| Item Formats and Scoring | 381 |
| Correction for Guessing | 381 |
| Item Distributions. | 382 |
| Relation of Items to the Scale as a Whole | 383 |
| Numerical Example | 387 |
| Test Reliability | 388 |
| The Logic of Internal Consistency Measures | 389 |
| | 2017 |

Introduction and Preview

Chapter Overview

Chapter 1 contains six major topics:

- I. MULTIVARIATE ANALYSIS: A BROAD DEFINITION—I begin by offering a very broad definition of multivariate analysis. The definition is: *an inquiry into the structure of interrelationships among multiple measures*.
- II. MULTIVARIATE ANALYSIS: A NARROW DEFINITION—The broad definition is important because virtually all behavioral research deals with questions about structure. Unfortunately, it would subsume all of statistical analysis, making for quite a large textbook! Consequently, it is customary to define multivariate analysis in more limited terms: the study of linear representations of relations among variables. Contained within this narrow definition are several interrelated models that are discussed in later chapters: (1) multiple regression (Chapters 4 and 5), (2) factor analysis (Chapters 6 and 7), (3) discriminant analysis, related classification techniques, and the multivariate analysis of variance and covariance (Chapters 8–10), (4) profile and canonical analysis (Chapter 11), and (5) analysis of scales (Chapter 12).
- III. SOME IMPORTANT THEMES—Models like factor analysis and multiple regression used to infer structure are all too often looked upon as if they were separate, unrelated procedures. In reality, certain themes recur regardless of the specific analytic procedure: (1) equations that are the byproduct of a particular analysis should be meaningful; (2) cutoffs need to be established defining "how high is high" and "how low is low" when one develops a prediction equation; (3) questions of statistical significance often arise; (4) outliers can render otherwise elegant analyses meaningless or, worse, highly misleading; (5) multivariate investigation should be guided by theory; and (6) particular problems are likely to arise when relations among individual items, as opposed to multi-item scales, are the unit of analysis.
- IV. THE ROLE OF COMPUTERS IN MULTIVARIATE ANALY-SIS—A generation ago, many analyses that were conceived of were not practicable. The availability of computers and sophisticated computer pack-

ages now make such analyses routine. The increasing role of the personal computer is also noted.

V. CHOOSING A COMPUTER PACKAGE—There are several important computer packages now on the market. Some of the pros and cons of the three most popular ones—SAS, SPSSX, and BMDP—are discussed.

VI. PROBLEMS IN THE USE OF COMPUTER PACKAGES—Knowing how to use a computer package is important to multivariate analysis, but it is not sufficient. Four necessary ingredients to a successful analysis are considered: (1) substantive knowledge of the research topic; (2) computer knowledge, including knowledge of the package; (3) empirical experience with various kinds of data; and (4) formal knowledge of the analytic procedures.

This textbook is written for those who need to analyze complex behavioral data, whether in field settings such as clinical and social psychology, in applied settings like nursing and marketing research, or in experimental settings like learning and perception. Although I cannot escape some degree of formality, I am writing for students whose primary interests lie more in empirical phenomena as opposed to rigorous mathematical statistics.

Although I very much appreciate why examples related to your specific interests are important, I hope I can "wean" you from examples limited to your own content area as the textbook topics progress. One of the things you should note is that if you are a clinical psychology student working with personality test data, you will have problems that are abstractly *identical* to a market researcher working on a consumer survey. One of the reasons that I find quantitative applications so rewarding is that I can work in a variety of areas, as long as I have someone to provide me with a background of the underlying empirical issue or have that background myself. At any one time, I can and have been working on the issue of selecting police officers, studying patients attitudes toward health care, looking at choices of financial institutions, and evaluating different forms of eye surgery. What I learned in one setting applied to all. I hope you can share the enjoyment that I have been experiencing.

One of the decisions that any author of a quantitative textbook needs to confront is how to choose examples. Conceptually, it should not matter if examples come from clinical psychology, marketing, or visual psychophysics; yet it does matter to most students. What I will do, therefore, is to pick my examples to illustrate the relevance of the various multivariate models in different settings.

Multivariate Analysis: A Broad Definition

As with most topics, it is useful to try to define what is meant by the term "multivariate analysis" before proceeding too far into its details. In the broadest and most literal sense, it means an inquiry into the structure of