

Cairo University
Faculty of Economics and Political Science
Department of Statistics

On Missing Values in Classification Analysis

310.9
M.A

By

64219

Mahmoud Al.Said Mahmoud



Under the Supervision of

Prof. Adel M. Zaher

Department of Statistics
Faculty of Economics
and Political Science
Cairo University

A.Prof. Ahmed H. Haroun

Department of Applied Statistics
Institute of Statistical
Studies and Research
Cairo University

A Thesis submitted to the Faculty of Economics and
Political Science (Department of Statistics)
in partial fulfillment of the requirement
for
The M.Sc. Degree in Statistics
1997



Acknowledgment

I am indebted to Prof. Adel Zaher for the considerable help , encouragement and valuable advice he has given me through his kind supervision of this dissertation. I would also like to express my gratitude to Dr. Ahmed Haroun for his effort and advice through the duration of this study. Their help and encouragement were particularly important to me and provided a constant source of inspiration.

Finally, I would like to give special thanks to my wife and my family for their support and prayers at all times.

Mahmoud Al Said.

Abstract

Classification Analysis is concerned with the problem of classifying a subject to one of several distinct groups on the basis of a set of measurements. For example, in business, a bank loan officer wishes to classify loan applicants to low risk credit customers or high risk credit customers on the basis of a set of variables. The most well-known and widely used rules of classification are presented through this thesis.

The presence of missing values in a data set used for building a classification rule is a serious problem that may face the investigator when applying the classification analysis (CA) to a practical situation. Many procedures have been developed to handle the missing values when applying (CA). The default method of handling missing values in (CA) used by many statistical packages (for example, SAS, Minitab and SPSS) is to omit all units containing missing values. Thus, considerable information may be lost due to the reduction of the sample size. The most commonly used procedures for this purpose are displayed in a separate chapter.

Many studies dealt with this problem in case of two multivariate populations with equal covariance matrices while a few studies treated it in case of two multivariate populations with unequal covariance matrices. The present study deals with the problem of classification analysis with missing values in case of two or more multivariate normal populations with equal and unequal covariance matrices through a simulation study. Three rules of classification and five methods of handling missing values are considered. The objective of this study is to compare the different

methods of handling missing values with respect to their ability in obtaining a “good” classification rule. In this thesis, two patterns of missing values are considered and the mechanism that lead to the presence of missing values is assumed to be missing at random (MAR).

Seven factors are taken into consideration. These factors are the sample size, the number of populations, the number of variables, the distance between the populations, the covariance matrices (equal or unequal), the pattern of missing values and the percentage of missing values. The impact of each factor on the methods of handling missing values is studied. A Minitab macro was designed to run the necessary calculations. This macro is displayed in a separate appendix.

Table of Contents

Chapter 1: Introduction	5
1.1 Review of Related Literature	6
1.2 Aims of the Thesis	9
1.3 Structure of the Thesis	10
Chapter 2: Classification Analysis	11
2.1 Rules of Classification	13
2.1.1 The Generalized Distance Rule (GDR)	13
2.1.2 The Quadratic Discriminant Function Rule (QDF)	14
2.1.3 The Linear Discriminant Function Rule (LDF)	16
2.1.4 The Maximum-Probability Rule	16
2.1.5 The Coefficient of Profile Similarity Rule (CPS)	17
2.1.6 The Logistic Regression Rule	19
2.1.7 Classification Trees	20
2.2 Comparison of Classification Rules.	22
2.3 Using Discriminant Analysis in Classification.	23
2.4 Estimating the Probability of Correct Classification.	24
2.5 Cost of Misclassification.	24
Chapter 3: Methods of Handling Missing Values	26
3.1 Patterns of Missing Values	27
3.2 Mechanisms of Missing Values	29
3.3 Procedures of Handling Missing Values	29
3.3.1 The Complete-Case (CC) Procedure	30
3.3.2 The Available-Data (AD) Procedure	31
3.3.3 The Estimation Procedure	31
3.3.3.1 The Mean Substitution Method (MS)	32
3.3.3.2 The Principal Component Estimation Method (PC)	32
3.3.3.3 The Regression Estimation Method (RG)	33
3.3.3.4 The Expectation Maximization Algorithm (EM)	34
Chapter 4: Methodology and Tools	38
4.1 The Two Populations Case	40
4.1.1 Equal Covariance Matrices	40
4.1.2 Unequal Covariance Matrices	41
4.2 The Three Populations Case	41

4.2.1 Equal Covariance Matrices	42
4.2.2 Unequal Covariance Matrices	43
4.3 Estimating the Probability of Correct Classification in the Complete Data Set.	43
4.4 Estimating the Probability of Correct Classification in the Incomplete Data Set.	44
Chapter 5: Results of The Simulation	47
5.1 Results of the Complete Data Set	48
5.2 Results of the Incomplete Data Set	55
Chapter 6: Summary and Conclusions	68
References	71
Appendix A: Tables of the Results of the Simulation	73
Appendix B: The Macro Program	102
Arabic Summary	



