# COMPARISON OF THE EUCLIDEAN AND LINEAR DISCRIMINANT FUNCTIONS IN STATISTICAL DISCRIMINANT ANALYSIS

A thesis presented to Massey University
in partial fulfillment of the requirements
for the degree of Master of Science
in Operations Research at
Massey University

Tiew-Kim Lim, BSc.

1992

# *Abstract*

*It is known that in the problem of statistical discriminant analysis, the linear discriminant function performs poorly when the dimension of the data, p, is large. It has been demonstrated by Marco, Young and Turner (1987) that the much simpler Euclidean distance classifier may out-perform the usual linear discriminant function under certain conditions. Their conclusions were arrived at from a simulation experiment which compared the probabilities of misclassification associated with the Euclidean distance classifier with those of the linear discriminant function, under certain conditions. In this dissertation, the asymptotic expansions of the probabilities of misclassification (the expected actual and expected plug-in error rates) associated with the two discriminant functions are obtained. These error rates are then used to investigate the relative performances of the two methods.*

*Chapter 1 introduces the problem of discriminant analysis and describes the two competing procedures for discriminant analysis and some associated error rates. Then Chapter 2 reviews previous results, in the literature which show that the Euclidean distance classifier can perform better than the linear discriminant function. Chapter 3 gives the asymptotic expansions of the error rates, i.e. the expected actual error rate, and the expected plug-in error rate. The relative performances of the two methods on the basis of the asymptotic expansions are discussed in Chapter 4. The results show that in general the plug-in error rates for the*

*Euclidean distance classifier give better estimates of the actual error rates for all dimensions of p which were considered, when compared to the linear discriminant function. Furthermore, the actual error rates for the Euclidean distance classifier also seem to give better estimates of the true error rates at large dimensions of p, when compared to the linear discriminant function. Certain situations where the linear discriminant function performs better than the Euclidean distance classifier are also identified. Final conclusions, discussions and recommendations for further work are given in Chapter 5.*

# *Acknowledgements*

First of all I would like to thank my supervisor, Dr. C.R.O. Lawoko for his supervision and encouragement, and for presenting the results of this dissertation at the International Symposium on Multivariate Analysis and its Applications in Hong Kong.

The department of Statistics at Massey University and the New Zealand Statistics Association sponsored my presentation of the results of this dissertation at the Annual Conference of the New Zealand Statistical Association at Victoria University of Wellington. This is hereby acknowledged.

I am also particularly grateful to Dr. S. Ganesalingam for his time and help looking through my algebra, and to Richard Rayner who helped with computing facilities.

A big thank you to Uncle K.C. and Aunty Lee Lee for their loan to make my studies possible. Thank you for your care, love, encouragement and prayers. Thanks also go to Yuin-Khai for loaning me his Smart Alec's Wally Jokes book to read when I was bored and frustrated with my work. I would also like to thank my parents and brother for their support and patience through my studies.

Finally, lots of thanks go to thank Simon, Jeanne and Ming for their friendship, care and encouragement.

# *Contents*

# List of figures and tables