

COMPARATIVE EVALUATION OF STORED-PATTERN
CLASSIFIERS FOR RADAR AIRCRAFT IDENTIFICATION

DISSERTATION

Presented in Partial Fulfillment of the Requirements for
the Degree Doctor of Philosophy in the Graduate
School of The Ohio State University

By

Sargur N. Srihari, B.Sc., B.E., M.S.

* * * * *

The Ohio State University

1976

Reading Committee:

Prof. L. J. White
Prof. A. A. Ksienski
Prof. S. H. Zweben

Approved by

A handwritten signature in cursive script, reading "Lee J. White". The signature is written in dark ink and is positioned above the printed name "Lee J. White".

Adviser

Department of Computer
and Information Science

Copyright by
Sargur N. Srihari
1976

To
My Parents

ABSTRACT

The purpose of this dissertation is to investigate aspects of design and evaluation of classification algorithms for a type of pattern recognition problem. The problem is one that is encountered in the design of a radar aircraft identification system. This problem differs from the usual situation encountered in statistical pattern recognition theory. Here classes are characterized by labelled training patterns that are considered to be ideal. The test pattern is assumed to be a noisy version of one of the ideal patterns, where noise is generated by a statistical process.

The situation where the noise perturbation is Gaussian is considered in detail. The Bayes classifier in this case has the form of the Sebestyen classifier. Techniques for efficient computer implementation, and evaluation of the Sebestyen classifier are investigated. A relationship between the Sebestyen classifier and a nearest-neighbor classifier that utilizes the ideal patterns as its reference set is obtained. When reference patterns satisfy certain location constraints, it is shown that the decision surfaces of the two classifiers are identical.

Classifier performance is evaluated in terms of misclassification probability, "classification bias", computational economy, and the effect of changes in classifier parameters. These parameters include dimensionality, assumed value of the noise level, and the number of subclasses in each class.

The problem of classifier design when more than one noisy pattern can be used to make a decision is investigated. This is the problem of efficient utilization of cumulative evidence to improve classifier performance. One of the procedures evaluated operates by taking a majority-vote over multiple decisions of a given classifier. It is shown that the misclassification probability of a majority-vote procedure increases monotonically with the "classification-bias" of the given classifier. Finally, the performance of such a procedure is compared to a procedure that utilizes the Sebestyen classifier in a sequential mode.

TABLE OF CONTENTS

	Page
ACKNOWLEDGEMENTS	v
VITA	vi
LIST OF FIGURES.	vii
LIST OF TABLES	ix
 Chapter	
I. INTRODUCTION AND OVERVIEW.	1
A. Problem Formulation.	2
B. Existing Work Related to this Research . .	6
C. Summary by Chapters.	10
II. THE RADAR AIRCRAFT IDENTIFICATION PROBLEM. . .	16
A. Pattern Representation Space	18
B. Pattern Dependency on Aircraft Orientation	23
C. Representation of Discrete Orientations. .	24
D. Classifiers for Radar Aircraft Identification	28
III. THE BAYES CLASSIFIER - ITS IMPLEMENTATION AND EVALUATION	30
A. The Bayes Classifier	31
B. The Bayes Classifier for the Radar Aircraft Identification Problem.	33
C. Computer Implementation of the Bayes Classifier	38
D. The Bayes Misclassification Probability. .	46

Chapter	Page
IV. THE DECISION SURFACES OF THE SEBESTYEN AND NEAREST-NEIGHBOR CLASSIFIERS	65
A. Nearest-Neighbor Classification.	66
B. Relationship Between the Bayes and Nearest-Neighbor Classifiers	68
C. Conditions for Similarity of the Sebestyen and Nearest-Neighbor Classifiers	77
1. Conditions for the Identity of the Classifiers.	77
2. Algorithm for Determining Identity . .	88
D. Robustness and Stability of the Sebestyen and Nearest-Neighbor Classifiers	92
E. Effect of Increasing Dimensionality on the Sebestyen and Nearest-Neighbor Classifiers	103
1. Dimensionality and the Sebestyen Classifier	106
2. Dimensionality and the Nearest- Neighbor Classifier.	111
V. CLASSIFICATION BIAS, MISCLASSIFICATION PROBABILITY, AND ECONOMY OF IMPLEMENTATION . .	120
A. Classification Bias.	121
B. Estimating the Misclassification Probability of the Nearest-Neighbor Classifier	126
C. Experimental Evaluation of the Nearest- Neighbor and Bayes Classifiers	130
1. Robustness of the Sebestyen Classifier	136
2. Effect of Increasing the Number of Subclasses	139
D. Economy of Implementation.	145
VI. CLASSIFICATION BASED ON MULTIPLE OBSERVATIONS..	150
A. An Estimate-Based Procedure.	152
B. A Suboptimal Procedure	153
1. Misclassification Probability of the Majority-Vote Algorithm.	156

Chapter	Page
C. Performance of Procedures with Radar Aircraft Data.	164
D. Sequential Classification.	167
1. Sequential Classification Using the Bayes Classifier	168
2. Application of Sequential Procedure to Radar Aircraft Identification . . .	175
VII. SUMMARY AND CONCLUSIONS.	179
A. Summary.	179
B. Specific Contributions and Conclusions . .	182
C. Future Research.	188
APPENDIX A	192
APPENDIX B	194
REFERENCES	198

ACKNOWLEDGMENTS

I would like to express my sincere gratitude to everyone who contributed to my graduate education and to the completion of this dissertation. Special thanks are due to the members of the reading committee: Professor A. A. Ksienski, Professor S. H. Zweben, and Professor L. J. White. In particular I would like to thank Professor Lee J. White for his guidance during the course of this work and for his untiring assistance in the preparation of this dissertation. I would like to acknowledge the support of this work on a grant from the Air Force Office of Scientific Research.

Professor B. Chandrasekaran deserves special appreciation for the inspiration he has given me and for his critical reading of the manuscript. I would like to thank Mrs. Judy Dettloff for patiently typing the many drafts of the dissertation. Finally my thanks are due to good friends Arthur Pyster, Coimbatore Satyapriya, and K. B. Lakshmanan for their assistance during the completion of this work.

VITA

- May 7, 1949 . . . Born - Bangalore, India
1967. B.Sc., Physics and Mathematics, Bangalore
University, Bangalore, India
1970. B.E., Electrical Communication Engineering,
Indian Institute of Science, Bangalore,
India
1971. M.S., Computer and Information Science, The
Ohio State University, Columbus, Ohio
- 1970-1973 Research and Teaching Assistant, Department
of Computer and Information Science, The
Ohio State University, Columbus, Ohio
- 1973-1975 Research Associate, Electrosience
Laboratory, The Ohio State University,
Columbus, Ohio
1975. Teaching Associate, Department of Computer
and Information Science, The Ohio State
University, Columbus, Ohio

PUBLICATION

"DABQUEST - A data base question answerer in PL/I", M.S.
Thesis, Department of Computer and Information Science, The
Ohio State University, December, 1971.

FIELDS OF STUDY

Major Field: Computer and Information Science

Studies in Pattern Recognition: Professors L.J. White,
A.A. Ksienski, B. Chandrasekaran
Studies in Computer Organization: Professor Ming-Tsan Liu
Studies in Information Storage and Retrieval: Professor
James E. Rush

LIST OF FIGURES

Figure		Page
Chapter II		
2-1	A radar aircraft identification system. .	17
2-2	Coordinate system to represent aircraft orientation	21
2-3	Functional representation of low-frequency radar return amplitude.	23
2-4	Typical data manifold in the radar aircraft identification problem	25
2-5	Discrete representation of orientation-space	26
Chapter III		
3-1	Behavior of Bhattacharyya coefficient with σ when $(\lambda_2/\lambda_1) < 1.86$	55
3-2	Behavior of derivatives of $f_1(\sigma^2)$ with σ^2	59
3-3	Behavior of $\frac{1}{2}\rho$ as a function of σ	62
Chapter IV		
4-1	Schematic illustration of classifiers C_{mult} and C_{mix}	70
4-2	Example of non-identity of C_{mult} and C_{mix}	72
4-3	Illustration of Sebestyen and nearest-neighbor decision surfaces.	75
4-4	Reference pattern configuration for obtaining conditions on classifier identity.	79
4-5	Examples of identity of the Sebestyen and nearest-neighbor classifiers.	82
4-6	Reference pattern configuration to illustrate dependence of the Sebestyen decision boundary on assumed σ	94

Figure		Page
Chapter IV (cont'd)		
4-7	Configuration of reference patterns to illustrate a negative effect of dimensionality.	108
4-8	Effect of dimensionality on the nearest-neighbor classifier	115
Chapter V		
5-1	Example of non-zero classification bias of the Bayes classifier	122
5-2	Evaluation of classifier performance with radar aircraft data: amplitude and phase.	131
5-3	Evaluation of classifier performance with radar aircraft data: amplitude only	132
5-4	Robustness of the Sebestyen classifier.	138
5-5	Performance of the Sebestyen classifier when class-conditional densities are Rician mixtures	140
5-6	Effect of increasing the number of sub-classes	142
5-7	Variation of nearest-neighbor decision surface with increased number of orientation possibilities	145
Chapter VI		
6-1	Variation of $p_i(n)$ with p_i	158
6-2	Functional behavior of $f(t)$	161
6-3	$P_e(n)$ as a function of n	163
6-4	Comparison of performance of fixed sample-size multiple observation procedures with radar aircraft data	165
6-5	The Bayes sequential procedure.	169

LIST OF TABLES

Table		Page
1.	Complexity of computing components of discriminant functions	146
2.	Summary of computational complexity of alternative classifiers.	147
3.	Performance of Bayes sequential procedure with radar data, $\sigma = 0.5$	177
4.	Performance of Bayes sequential procedure with radar data, $\sigma = 1.0$	177

CHAPTER I

INTRODUCTION AND OVERVIEW

The purpose of this dissertation is to investigate aspects of design, implementation, and evaluation of classification algorithms for the pattern recognition problem encountered in radar aircraft identification. The classification algorithms are based on the approach of feature extraction followed by statistical classification. When this approach is taken, a pattern is a vector of feature measurements, and classification is accomplished by partitioning the feature-space. The feature-space partitioning is accomplished by utilizing a set of training patterns that are representative of the different categories. The classifier associates a category label with each of the regions of the partitioned feature space. The pattern to be classified, called the test pattern, is assigned by the classifier to that category whose label is attached to the region of feature-space in which the test pattern lies.

In most pattern recognition problems that are solved by taking the approach of feature extraction and statistical classification, classes are described by statistical distributions. When training patterns are labelled, the usual method of designing a classifier is to determine the parameters of class distributions and to use these parameters in determining the discriminant functions of the classifier. An alternative method of classifier design is that of storing individual training patterns as reference patterns, and performing classification by comparing the test pattern to each of the reference patterns. The latter method works well with many intuitively appealing rules of comparison, as long as a large and representative reference pattern set is available. However, it tends to be wasteful of both storage space and time taken to make decisions, particularly when the parameters can be a priori estimated.

I.A. Problem Formulation

The problem considered in this research differs from the usual pattern recognition problem in a number of ways. In this problem classes are characterized by labelled

training patterns that are considered to be "ideal". The test pattern is assumed to be a "noisy" version of one of the "ideal" patterns, where noise is generated by a statistical process. Further, it is assumed that the "ideal" patterns are generated by a set of functions. Each category is characterized by a different function, and every function is defined over the same domain. In the radar aircraft identification problem the function domain is the set of all possible values of the orientation coordinates of an aircraft with respect to the radar. At the time of classification the aircraft orientation is assumed to be indeterminate. In other words, values of the orientation coordinates are either unknown or determined only approximately. Essentially the classification problem is one of composite hypothesis testing for a parameter that represents the "ideal" pattern.

One of the objectives of this research concerns the implementation and evaluation of the Bayes classifier that is based on a decision-theoretic model of the problem. In particular, the situation where noise perturbation of the ideal pattern is Gaussian will be considered at length.

The implementation part of the problem is to obtain an algorithm for efficient computation and comparison of Bayes discriminant functions. The evaluation part is primarily to obtain analytical bounds on the misclassification probability of the Bayes classifier, and to characterize the behavior of misclassification probability with the value of the noise parameter. Another aspect of classifier evaluation is to examine the effect of the assumed value of the noise parameter on classifier performance.

A second objective of this research is to relate the Bayes classifier for the problem to a family of classifiers based on the nearest-neighbor rule. This includes comparison of decision surfaces of the classifiers, and obtaining conditions on the configuration of reference patterns for their misclassification probabilities to be nearly the same. The classifiers will also be evaluated in terms of other criteria that are relevant to the type of problem being considered. These criteria include average misclassification probability, computational economy in typical implementations, and the effect of changes in classifier parameters. The classifier parameters to be considered

are: dimensionality, standard deviation of Gaussian noise and uncertainty associated with the measurement of orientation. In particular, the effect of increase in dimensionality will be examined from the point of view of characterizing situations that lead to increased misclassification probability at particular orientations.

The third objective of this research is to examine the problem of classifier design when a sequence of noisy patterns can be used to make a decision. In this problem, each of the observed patterns can be noisy versions of reference patterns corresponding to different orientations. As a consequence more complex procedures than simple hypothesis testing will have to be used. Among different procedures to be considered is one that takes a majority-vote over multiple decisions of a given classifier. The misclassification probability of such a procedure is dependent on the misclassification probability of the classifier with which it is implemented.

I.B. Existing Work Related to this Research

Most of the results obtained and techniques used in this research have been based upon a considerable body of work done by other researchers. In this section we will discuss such work and indicate its relationship to this dissertation. The relevant literature can be grouped into two categories. The first category is the literature pertaining to the problem of classifier design for radar identification of targets and aircraft. The second category is the large body of work concerning classification theory.

The problem of classifier design for radar identification of objects of arbitrary shape, referred to as targets, has received some attention in the recent past. Repjar (1970) obtained a linear discriminant classifier which separates classes by hyperplanes in feature-space. The hyperplanes were obtained such that the misclassification probability between classes is minimized when reference patterns are corrupted by Gaussian noise. Considering the aircraft identification problem, Lin (1974) showed that the linear discriminant classifier performs well only when the region of uncertainty in the orientation-space is small. The

reason for the poor performance of the linear discriminant classifier is the highly convoluted nature of the data surfaces due to which separating hyperplanes cannot be found. Even when it is possible to find separating hyperplanes, the misclassification probability between classes tends to be considerably high. Further, Lin demonstrated that significant improvement in performance can be obtained by utilizing the nearest-neighbor classifier in a situation where noise is Gaussian. White and Ksienski (1974) showed how a data surface approximation that passes through a set of reference patterns can be utilized in conjunction with the nearest-neighbor rule. They also demonstrated that when the data surfaces satisfy certain conditions, one of the classes is more susceptible to misclassification than others.

The Bayes classifier examined in this research closely resembles the classifier investigated by Sebestyen (1962) in the design of automated systems for speaker verification. In the problem of speaker verification each utterance of a speaker corresponds to a different region of the relevant feature-space. Here the sounds uttered by a given speaker lie on a multiply-connected domain with

considerable overlap between the domains of two speakers. The classifier for the problem was implemented by determining the necessary parameters, viz., "ideal" patterns, and variances of component Gaussian distributions, by an adaptive process.

The problem of evaluation of misclassification probability of classifiers has received a great deal of attention recently. Toussaint (1974) provides a bibliography of nearly two hundred publications, most of which pertain to the problem of having few training patterns. Bounds on the misclassification probability of classifiers in terms of various statistical distance measures are available in the literature (see Chen, 1973). In particular, the bound obtained by Kailath (1967) in terms of the Bhattacharyya coefficient between statistical distributions is utilized in this research. This bound is used to obtain a simple bound on the misclassification probability of the Sebestyen classifier in terms of the maximum and minimum inter-class separation of patterns.

The problem of comparing classifier decision surfaces has received considerably less attention than the problem of estimation of misclassification probability. This is

perhaps due the utility of such a comparison being primarily in understanding classifier performance. Batchelor (1973) compared the decision surfaces of the nearest-neighbor and potential function classifiers. He obtained sufficient conditions on reference pattern location for which the two decision surfaces are similar.

The problem of robustness of a classifier with non-normality has been examined by Lachenbruch, et. al. (1973). They showed that the linear discriminant function classifier is greatly affected by certain types of non-normality. Due to its complexity, the Sebestyen classifier does not seem to have been examined from the point of view of robustness with non-normality.

One of the criteria of performance evaluation of classifiers is the computational effort needed to make decisions. Economy of implementation of various classifiers have been studied by Hahn and Nelson (1972). They compared the economy of implementation of certain parametric classifiers based on Gaussian and Rayleigh distributions. However, evaluation of relative computational economy in using the Sebestyen classifier over the nearest-neighbor classifier does not appear to have been made.

The problem of classifier design with multiple observations was considered by Sebestyen (1962). He showed how the misclassification probability of a classifier based on a likelihood ratio test can be expressed in terms of the misclassification probability of the same classifier based on the use of a single observation.

Considering the problem of utilizing as few observations as necessary, Roberts and Mullis (1970) demonstrated how a single observation Bayes classifier can be utilized in the design of a sequential procedure.

I.C. Summary by Chapters

Throughout the discussion the radar aircraft identification problem is used as the paradigm. In order to acquaint the reader with the problem of radar aircraft identification, a hypothetical system is described in Chapter II. This includes a discussion of the data structure of the problem that makes it necessary to consider more complex classification procedures than linear discriminant techniques.

In Chapter III, the classification problem is formulated in terms of Bayes decision theory. The formulation

includes a statistical model of uncertainty in the values of variables in the function domain. The Bayes classifier based on the decision-theoretic model has the form of the Sebestyen classifier under certain conditions. The discriminant functions of the Sebestyen classifier are sums of exponentials, which leads to a problem of computer implementation. A discriminant comparison algorithm that makes use of data values in order to avoid computing exponentials is discussed. The latter part of the chapter pertains to the problem of determining the Bayes misclassification probability. This includes obtaining bounds on the misclassification probability of the Sebestyen classifier in terms of the inter-class separation of patterns. Another aspect is that of predicting the dependency of misclassification probability of the Sebestyen classifier on noise standard deviation.

The decision surfaces of the Sebestyen and nearest-neighbor classifiers are examined in detail in Chapter IV. Conditions on the relative location of reference patterns for the similarity of performance of the two classifiers are discussed. The implementation of the Sebestyen classifier depends upon prior knowledge of the standard

deviation of noise. The dependence of the decision boundary on the assumed value of this parameter is examined analytically in the univariate case. It is shown that an underestimation of noise standard deviation will lead to a smaller change from the optimal decision boundary than an overestimation. The last part of Chapter IV deals with the issue of dimensionality and its effects on classifier behavior. The effect of increasing the dimensionality of either the Sebestyen or nearest-neighbor classifier is examined from the point of view of the resulting changes in decision surfaces and misclassification probabilities. Due to the presence of multiple subclasses certain counter-intuitive phenomena occur with change in dimensionality. The misclassification probabilities at particular orientations are seen to increase with an increase in the number of feature measurements. Sufficient conditions for such a phenomenon to occur are discussed. These conditions are on reference pattern configuration and on the discriminating capacity of the newly introduced feature.

In Chapter V a new criterion for classifier evaluation, termed classification bias, is defined. This criterion is

closely related to the average misclassification probability of a given classifier. An experimental evaluation of the misclassification probability and classification bias of the Sebestyen and nearest-neighbor classifiers with radar aircraft data is made. The average misclassification probabilities of the two classifiers are seen to be nearly the same with radar aircraft data. However, the Sebestyen classifier tends to have a lower classification bias than the nearest-neighbor classifier.

In the second part of Chapter V, experimental results that demonstrate the behavior of misclassification probability with changes in classifier parameters are reported. One set of experimental results pertain to the effect of incorrect assumption of subclass standard deviation on the average misclassification probability of the Sebestyen classifier. This set of results, obtained in a multivariate situation, reinforce the conclusions made with reference to the univariate situation in Chapter IV. Another set of experimental results provide insight into classifier behavior with changes in yet another parameter, viz., degree of uncertainty associated with orientation measurements.

Here again certain counter-intuitive phenomena are observed. The last part of Chapter V discusses the economy of implementation of different classifiers. It is observed that the computational economy in using the nearest-neighbor classifier over the Sebestyen classifier tends to be marginal as feature-space dimensionality is increased.

The subject of Chapter VI is the design of efficient procedures for the multiple observation problem. Among procedures examined is one that operates by taking a majority-vote over n individual decisions of a given classifier. Analytical expressions for the misclassification probabilities of such a procedure for the two-class and multi-class problems are obtained. It is shown that the misclassification probability of a majority-vote procedure is lowest if for a specified value of misclassification probability of the classifier, its classification bias is minimized. If the orientation remains constant as each of n observations are made, it is demonstrated that a procedure which submits an n -sample estimate to a given classifier performs better than a procedure which takes a majority-vote over n individual decisions of the classifier.

The last part of the chapter examines the design of sequential procedures for the multiple observation problem. An implementation of the Sebestyen classifier in a sequential mode is discussed. The stopping boundaries of such a procedure are obtained from standard sequential hypothesis testing procedures. The sequential procedure is conservative in that the experimental misclassification probabilities tend to be lower than the theoretical misclassification probabilities used to compute the stopping boundaries.

Chapter VII summarizes the major conclusions to be drawn from this research. Areas of future research that pertain to classifier design for highly convoluted ideal data surfaces are indicated.

CHAPTER II

THE RADAR AIRCRAFT IDENTIFICATION PROBLEM

The problem of radar aircraft identification (RAI) may be characterized by an operational system as shown in Fig. 2-1. The objective of the system is to determine the aircraft type under the following conditions. It is assumed that the aircraft to be identified is tracked by two different modes of radar. The first of these is a low-frequency radar, by which the amplitude and phase of radar return signals at a number of different low frequencies are measured. This data, expressed as a vector, constitute the pattern vector. Another input to the radar aircraft identification system is a measurement of the orientation of the aircraft with respect to the radar tracking station. The orientation measurement is obtained by conventional high-frequency radar, the details of which we will not be concerned with, except to note that it depends upon the track record of the aircraft (range and Doppler) and

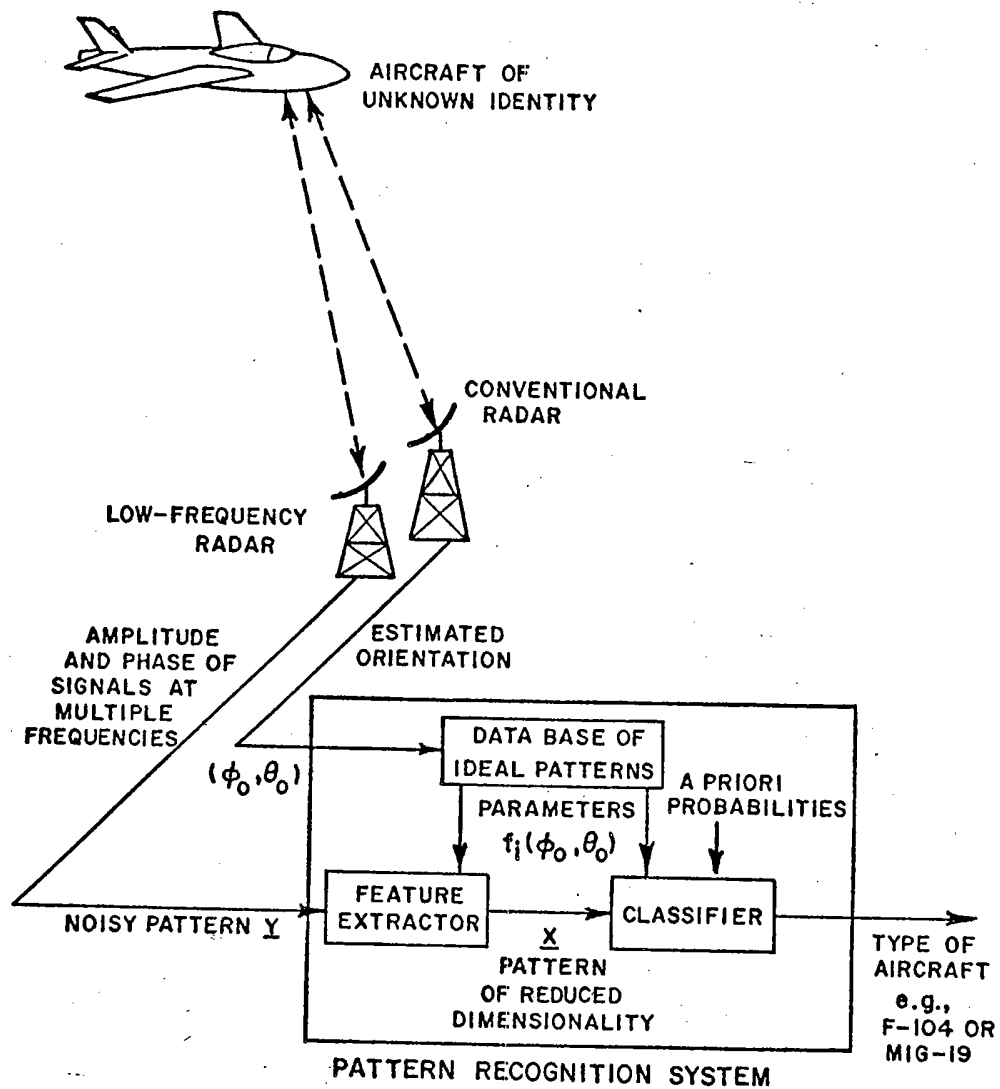


Fig. 2-1. A radar aircraft identification system.

flight dynamics.

The function of the feature extractor in Fig. 2-1 is to select the fewest number of features that enable good class discrimination at the estimated orientation. Feature extraction may be an on-line process in which case the objective is dimensionality reduction, for increased processing speed. This objective can be accomplished either by choosing a subset of the given measurements, or by performing a transformation between the original feature-space and a space of lower dimensionality. Feature extraction can also be an off-line process that is performed prior to the time of identification. The primary objective of off-line feature extraction is to select a specified number of frequencies out of a larger set of possibilities so as to satisfy radar hardware constraints. In either case the feature extractor obtains features on the basis of separability of classes in the resulting feature space.

The data-base of ideal patterns consists of noise-free values of the amplitude and phase of low-frequency radar returns at discrete orientations of each of the aircraft types considered. These values constitute the parameters, or reference set, of the classifier. The values can be

obtained by setting-up scaled model aircraft and measuring the returns at appropriately scaled frequencies. Alternatively, they can be computed by electromagnetic considerations. One such procedure for computing radar scattering data is by the use of wire-grid models of aircraft (Lin and Richmond, 1975).

An important part of designing any pattern recognition system is that of pattern analysis. This consists of subjecting the data to a variety of procedures for inferring structures that are present in the data. Procedures for such exploratory data analysis include the use of histogram plots, scatter plots, cluster analysis routines, data display routines, and linear discriminant analysis. The purpose is to understand the regularities and peculiarities of the data base to enable better classification. The rest of this chapter is devoted to the pattern analysis of the radar aircraft identification system.

II.A. Pattern Representation Space

The amplitude and phase of low-frequency radar return signals from the aircraft to be identified are assumed to be available at each of d frequencies. This set of values,

expressed as a column vector, is considered to be a pattern vector or simply a pattern. If amplitudes alone are used as features for class discrimination, then each of the elements of the pattern vector are real-valued and positive. If both amplitude and phase are used, then the elements are complex-valued. However, if the real and imaginary parts of each of the complex-valued elements are considered to be distinct components, then the pattern can be expressed as a 2d-dimensional vector with only real components. The effective dimensionality of the feature space is therefore d or $2d$ depending on whether amplitude alone, or amplitude and phase are used to construct the pattern vector.

The measured value of a pattern due to an aircraft depends upon the orientation of the aircraft with respect to the radar. In order to uniquely describe this orientation, a plane-oriented spherical coordinate system is used. In this coordinate system, as shown in Fig. 2-2, the origin is at the center of the aircraft, the nose aligned along the z -axis and the right wing along the x -axis. The direction of the radar relative to the aircraft is uniquely described by two angles ϕ and θ . The pair (ϕ, θ) is called the aspect angle and the vector space with coordinates ϕ and θ is the

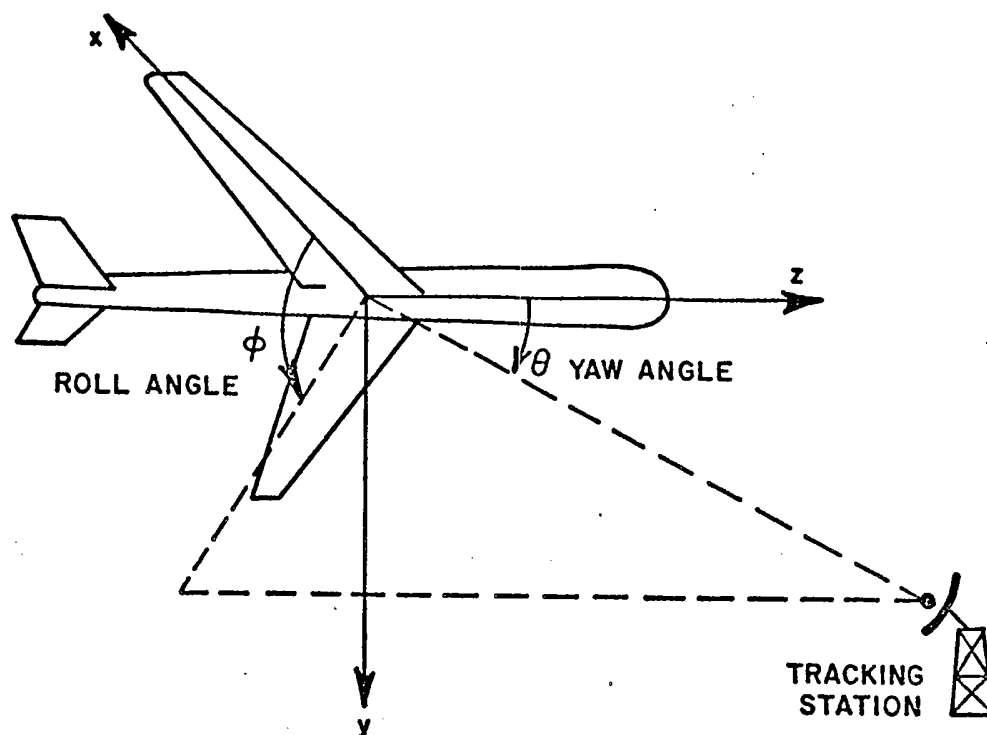


Fig. 2-2. Co-ordinate system to represent object orientation.

orientation space or the aspect angle space. Since the aircraft are symmetric about the y-z plane, the ranges defined by $-90^\circ \leq \phi \leq 90^\circ$, $0^\circ \leq \theta \leq 180^\circ$ specify all possible orientations.

Data used to obtain empirical results in this research were those that were theoretically obtained using wire-grid models of aircraft. The data correspond to returned signals from aircraft when vertically polarized low-frequency waves are used; the low frequencies considered were in the range from 2 to

24 MHz, as opposed to conventional radar frequencies that lie in the range from 10^3 to 30×10^3 MHz. Data for four aircraft types were available, these being the F-104, F-4, MIG-19, and MIG-21. For a listing, display, and further discussion of the data used in this research see Srihari and White (1975).

II.B. Pattern Dependency on Aircraft Orientation

The amplitude and phase of low-frequency radar returns from an aircraft at a given frequency may be regarded as a complex-valued function of the orientation of the aircraft. A graph of the dependence of ideal, i.e., noise-free, radar return amplitude on aircraft type and orientation is given in Fig. 2-3. The curves shown are for four aircraft types, with a fixed value of θ , and at a single frequency.

In Fig. 2-3 observe that the amplitude return for the MIG-19 in the aspect angle range of $\theta=120^\circ$, $60^\circ \leq \phi \leq 90^\circ$ is approximately twice as large as the corresponding return for the other aircraft. Therefore, in this aspect angle range, a single frequency is sufficient to distinguish the MIG-19 from other aircraft provided the noise level is not very high. Similarly, in the aspect angle range of

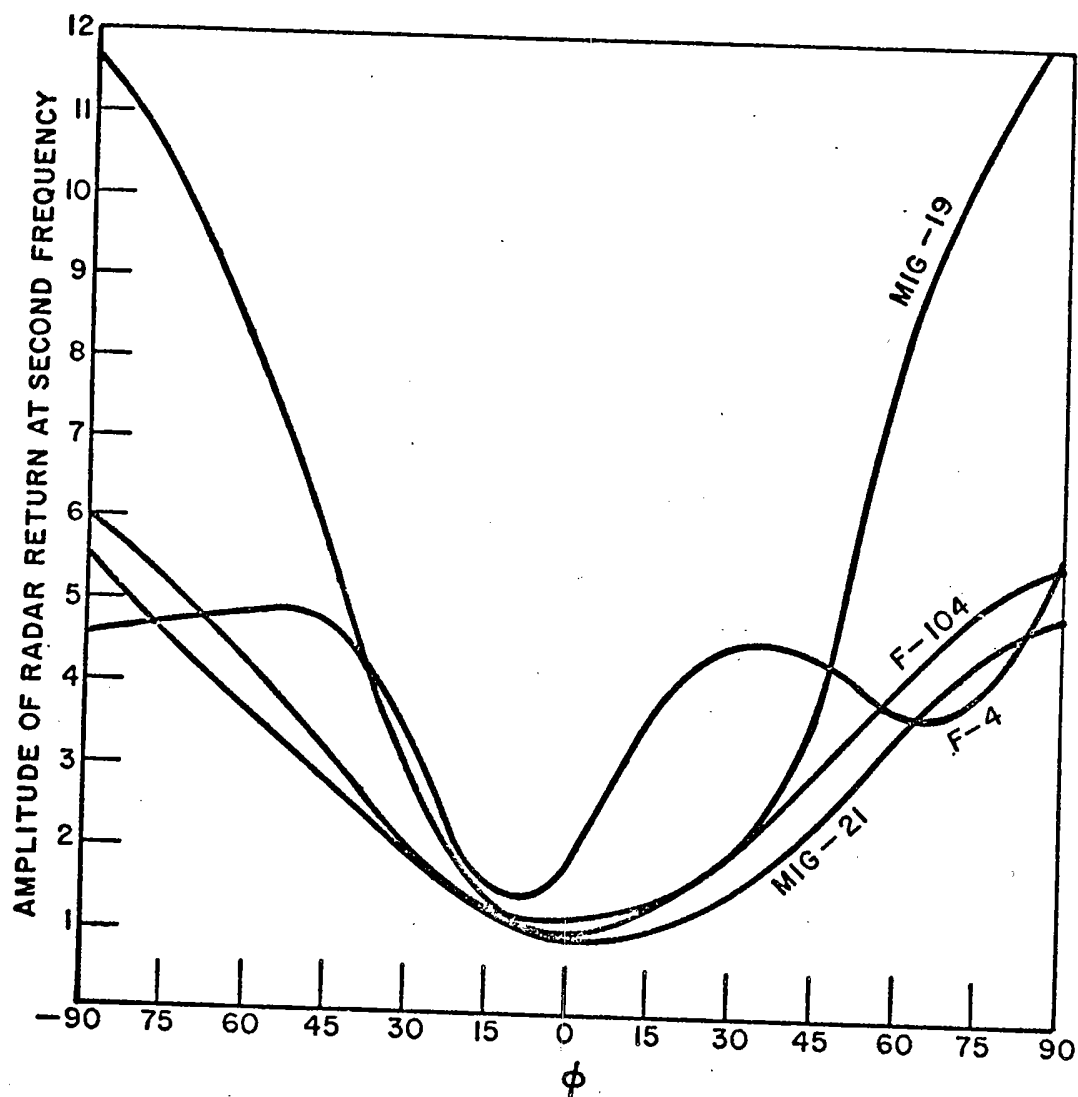


Fig. 2-3. Functional representation of low-frequency radar return amplitude.

$\theta=120^\circ$, $15^\circ \leq \phi \leq 30^\circ$, the F-4 is easy to discern. However, all other aircraft at aspect angles other than those in the above ranges pose a classification problem that requires more data.

A situation where two aircraft give rise to different signal returns at the same orientation while providing identical returns at different orientations can be seen in Fig. 2-3. For example, the F-4 at $(30^\circ, 120^\circ)$ is clearly distinct from the F-104 at $(30^\circ, 120^\circ)$, but is indistinguishable from the F-104 at $(66^\circ, 120^\circ)$. This poses a data separation problem which can be resolved by using orientation information; alternatively more features can be used to separate the data.

II.C. Representation of Discrete Orientations

Since there are two parameters underlying the data, viz., ϕ and θ , the ideal pattern vector for a given class must lie on a two-dimensional manifold imbedded in feature-space. An illustration of such a manifold, also called a data surface, is given in Fig. 2-4 for $d=2$. The data surface for class ω_i is characterized by the set of functions (f_{i1}, \dots, f_{id}) . With aircraft data these functions are

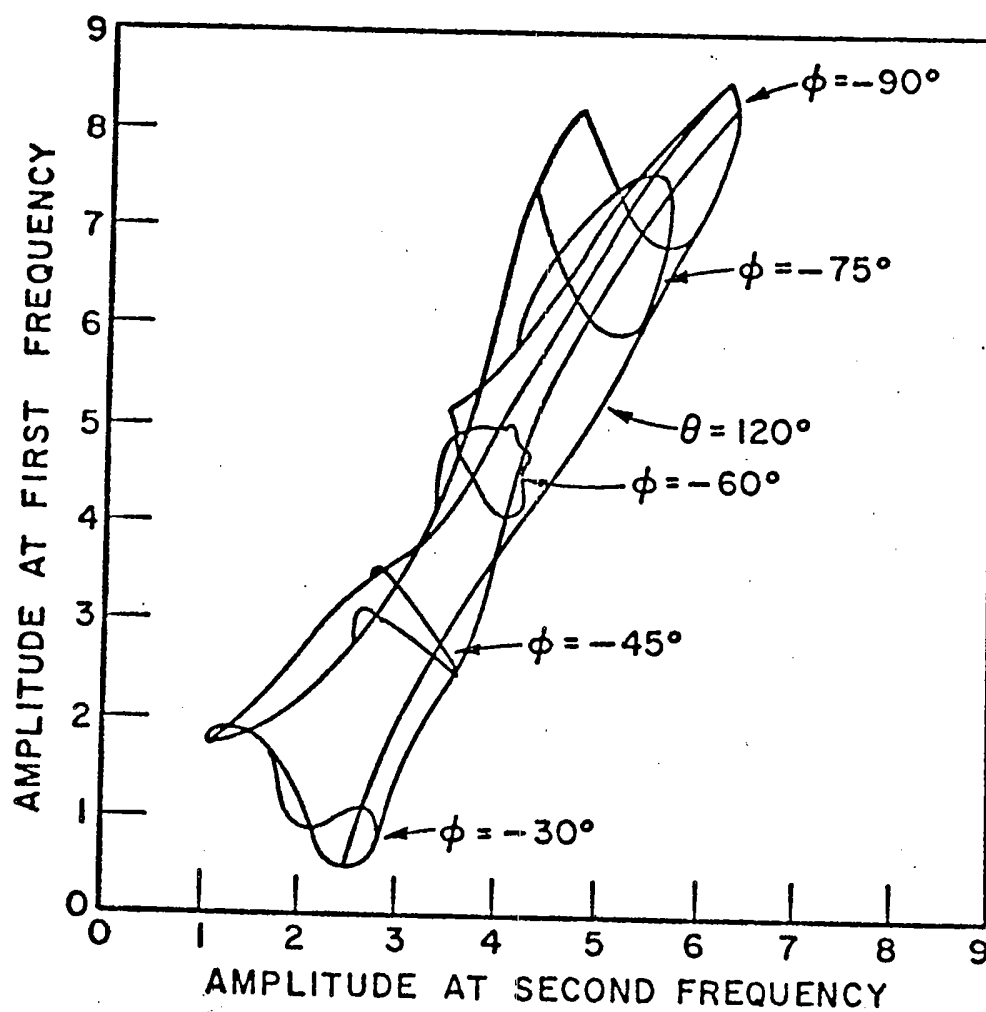


Fig. 2-4. Typical data manifold in the radar aircraft identification problem; data from the F-104 model, ϕ varying from -90° to -30° , θ varying from 0° to 180° .

not topological mappings, as observed from the fact that the data surface passes through itself. That is, every ideal pattern of a given class does not have a unique corresponding point in the domain of possible aircraft orientations. In the radar aircraft identification problem, the analytical forms of the functions are not known. However, the values of the functions at discrete values of (ϕ, θ) are known. In order to represent discrete points in orientation-space the following notation is introduced.

The orientation-space is considered to be a two-dimensional grid of $n_1 \times n_2$ points with ϕ ranging from $-\pi/2$ to $+\pi/2$ and θ ranging from 0 to π (Fig. 2-5).

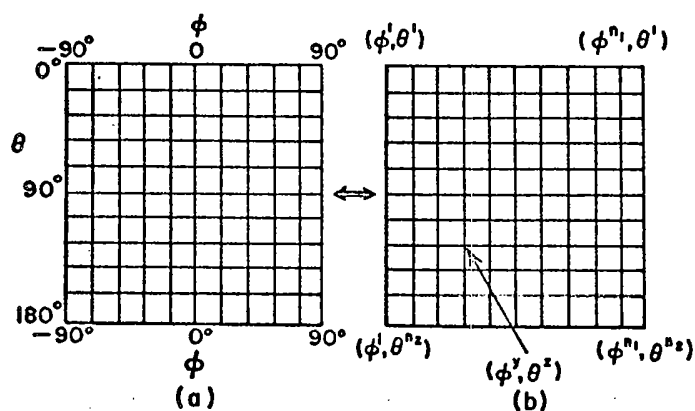


Fig. 2-5. Discrete representation of orientation-space.

In this notation the points are labelled as (ϕ^y, θ^z) , $y=1, \dots, n_1$; $z=1, \dots, n_2$, as shown in Fig. 2-5. The labels y and z are obtained from the actual values of ϕ and θ by the mapping functions

$$(2-1) \quad y = 1 + \frac{(n_1-1) \cdot [(\pi/2) + \phi]}{\pi} \quad \text{and} \quad z = 1 + \frac{(n_2-1)\theta}{\pi}.$$

If (ϕ_0, θ_0) is the observed value of the aircraft orientation then the true aircraft orientation, (ϕ, θ) , is assumed to be distributed according to the probability density function $q(\phi, \theta / \phi_0, \theta_0)$. In earlier implementation of classifiers (see Lin, 1974), this distribution has been considered to be uniform in the range $\phi_0 \pm 15^\circ$, and $\theta_0 \pm 10^\circ$. In general terms the discrete probability density function representing measurement uncertainty in orientation knowledge can be written as

$$(2-2) \quad q(\phi, \theta / \phi_0, \theta_0) = \frac{1}{(2a+1)(2b+1)} \quad \begin{array}{l} \phi = \phi^{y-a}, \dots, \phi^a, \dots, \phi^{y+a} \\ \theta = \theta^{z-b}, \dots, \theta^z, \dots, \theta^{z+b} \end{array}$$

where y and z are obtained from definition (2-1) with $\phi = \phi_0$ and $\theta = \theta_0$, and a, b are integers whose values depend upon n_1, n_2 , and the resolution of orientation measurement.

II.D. Classifiers for Radar Aircraft Identification

In order to select classification procedures for the problem of radar aircraft identification, parametric and non-parametric classifiers have been implemented and evaluated with low-frequency radar data by other researchers. These include a classifier based on linear discriminant functions (Repjar, 1970), and a classifier based on the nearest-neighbor rule (Lin, 1974). The linear discriminant function classifier for the problem was obtained by determining hyperplanes that separate classes such that the misclassification probability is minimized. This classifier was found to perform well when the number of discrete orientations considered is relatively small. When the number of orientation possibilities was increased, it was found that the misclassification probability increased considerably. The latter fact is explained by the highly convoluted nature of the data surfaces, due to which hyperplanes that separate the ideal patterns of classes cannot be found.

The classifier based on the nearest-neighbor rule, as applied to the radar aircraft identification problem, operates by measuring the Euclidean distance between the measured noisy pattern and each of the stored reference

patterns. At the time of classification the reference patterns considered by the classifier are those that correspond to orientations that have a non-zero probability of being the true orientation. The decision of the classifier is taken to be that class which has a reference pattern nearest, in terms of the measured distance, to the noisy pattern. This classifier was found to have a significantly lower average misclassification probability than the classifier based on the linear discriminant function (Lin, 1974). The effective decision surface of the classifier based on the nearest-neighbor rule is more complex than that of the linear discriminant classifier. Properties of the decision surface of the nearest-neighbor classifier will be discussed in Chapter IV.

CHAPTER III

THE BAYES CLASSIFIER - ITS IMPLEMENTATION AND EVALUATION

In this chapter the classification problem is formulated in terms of Bayes decision theory. A classifier which minimizes the criterion of average misclassification probability, called the Bayes classifier, is obtained. The Bayes classifier is expressed in terms of discriminant functions with two arguments, viz., the input pattern vector and the orientation measurement. The discriminant functions of the Bayes classifier have the same analytical form as those of the Sebestyen classifier when the noise perturbation is Gaussian. An algorithm for efficient computer implementation of the Sebestyen classifier is examined. This algorithm can also be used to implement other Bayes classifiers based on class-conditional distributions that are exponential mixtures. The latter part of the chapter concerns the misclassification

probability of the Sebestyen classifier. The dependence of misclassification probability on the interclass separation of ideal patterns and on the standard deviation of Gaussian noise is examined.

III.A. The Bayes Classifier

It is assumed that the pattern to be classified can belong to one of c predetermined classes labeled $\omega_1, \dots, \omega_c$. The other assumptions of the model are as follows.

- a. The d feature measurements constituting the ideal pattern vector of a given class are assumed to be functions over a domain \mathcal{J} . The domain consists of the set of all possible aircraft orientations with respect to the radar. The ideal pattern vector corresponding to class ω_i , and orientation $(\phi, \theta) \in \mathcal{J}$ is written as

$$\underline{f}_i(\phi, \theta) = [f_{i1}(\phi, \theta), \dots, f_{id}(\phi, \theta)] ,$$

where $f_{ij}(\phi, \theta)$ is the ideal feature measurement for the j^{th} feature.

- b. If (ϕ, θ) is the true orientation and ω_i the true pattern classification, then the distribution of the pattern that has been corrupted by noise, i.e., the measured pattern x , is given by the conditional probability density function $p(\underline{x}/\omega_i, \phi, \theta)$. The parameter vector of this density function is $P = [\underline{f}_i(\phi, \theta), \sigma^2]$ where σ^2 is the variance of additive Gaussian noise that perturbs radar return measurements.
- c. If (ϕ_o, θ_o) is the observed value of the orientation, then the true orientation has a distribution defined by the probability density function $q(\phi, \theta/\phi_o, \theta_o)$. This distribution represents the uncertainty in knowledge of orientation due to measurement or quantization error.
- d. The a priori probability of class ω_i is denoted by $P(\omega_i)$.

The average misclassification probability of a classifier is the sum of the weighted probabilities of not choosing the correct class, where weights are class a priori

probabilities and the summation is over classes. The classifier which minimizes the average misclassification probability, viz., the Bayes classifier, can be shown to have discriminant functions of the form

$$(3-1) \quad g_i(\underline{x}, \phi_o, \theta_o) = P(\omega_i) \sum_j p(\underline{x}/\omega_i, \phi, \theta) q(\phi, \theta/\phi_o, \theta_o),$$

where the summation is over all values of (ϕ, θ) at which $\underline{f}_i(\phi, \theta)$ are known. Eq. (3-1) follows from the fact that the Bayes classifier is one which chooses that class whose a posteriori probability is maximum (see Duda and Hart, 1973). The Bayes classifier is one that effectively chooses that class for which Eq. (3-1) is maximum.

III.B. The Bayes Classifier for the Radar Aircraft Identification Problem

In the radar aircraft identification problem appropriate functional forms for the probability density functions, p and q , can be determined and the Bayes discriminant functions evaluated. We have already seen that uncertainty in knowledge of the true orientation can be modeled by a probability distribution in the orientation-space. This is the distribution defined by Eq. (2-2). We will consider here the functional form of the probability density function of the noisy pattern.

Noise that affects radar return signals is generally assumed to be narrow-band Gaussian (Whalen, 1971). Thus the signal corresponding to a particular center frequency w_j is of the form

$$\begin{aligned} f(t) &= [A \cdot \cos(w_j t + p)] + v(t) \\ &= [A \cdot \cos p + v_1(t)] \cdot \cos(w_j t) - [A \cdot \sin p + v_2(t)] \cdot \sin(w_j t) \end{aligned}$$

where $v(t)$ has quadrature components $v_1(t)$ and $v_2(t)$ which are Gaussian processes with zero mean and variance σ^2 ; A and p are the uncorrupted values of the amplitude and phase of the signal. Denote

$$x_c(t) = A \cdot \cos p + v_1(t)$$

$$x_s(t) = A \cdot \sin p + v_2(t) .$$

Dropping explicit time dependence, x_c and x_s may each be regarded to be the noisy values of the returned signal along orthogonal coordinates at frequency w_j , where

$$v_1, v_2 \sim n(0, \sigma^2) .$$

To discriminate between classes on the basis of the returned signal at frequency w_j , a classifier which uses both amplitude (A) and phase (p) can regard $A \cdot \cos p$ and $A \cdot \sin p$ as the feature values along orthogonal coordinates. Since every frequency contributes two coordinates to the

feature-space, a classifier based on the complex signal returns at d frequencies involves a feature-space of dimensionality $2d$. Since noise is additive Gaussian along each dimension, the conditional probability density function of the noisy pattern at a given orientation is multivariate Gaussian. Assuming the effect of noise on the returns at each frequency to be independent,

$$(3-2) \quad p(\underline{x}/\omega_i, \phi, \theta) = \frac{1}{(2\pi)^d \sigma^{2d}} \exp \left\{ -(1/2\sigma^2) \left[\|\underline{x} - \underline{f}_i(\phi, \theta)\|^2 \right] \right\}.$$

Substituting Eqs. (2-2) and (3-2) into Eq. (3-1), the Bayes discriminant function takes the form

$$(3-3) \quad g_i(\underline{x}, \phi_0, \theta_0) = P(\omega_i) \sum_{\phi^Y=a}^{\phi^Y+b} \sum_{\theta^Z=b}^{\theta^Z+b} \exp \left\{ -(1/2\sigma^2) \left[\|\underline{x} - \underline{f}_i(\phi, \theta)\|^2 \right] \right\}.$$

A classification algorithm that uses the general form of (3-3) as its discriminant functions, i.e.,

$$g_i(\underline{y}) = \sum_{j=1}^S M_j \exp \left[-Q_j \|\underline{y} - \underline{a}_{ij}\|^2 \right]$$

where \underline{a}_{ij} are the reference patterns of ω_i , Q_j and M_j are positive scalars, and \underline{y} is the measured pattern, has been referred to as the Sebestyen Classifier (SC) in the literature (Batchelor, 1974).

Thus the Bayes classifier for the radar aircraft

identification problem, when amplitude and phase are used to construct the pattern vector, is the Sebestyen classifier with $M_j=1$ and $Q_j=(1/2\sigma^2)$. The Sebestyen classifier can be employed even if noise is not Gaussian, as long as training patterns are available, and reasonable assumptions for M_j and Q_j can be made. For example, if the ideal patterns \underline{a}_{ij} are not known, but the subclass associated with \underline{a}_{ij} is characterized by a set of k_{ij} noisy reference patterns, then the Sebestyen classifier can be implemented as follows: $M_j=k_{ij}$, $Q_j=(1/2\hat{\sigma}_{ij}^2)$, $\underline{a}_{ij}=\hat{\underline{a}}_{ij}$, where $\hat{\sigma}_{ij}^2$ and $\hat{\underline{a}}_{ij}$ are the sample variance and sample mean of the k_{ij} patterns respectively. This approach to the implementation of the Sebestyen classifier has been shown to be useful in the problem of speaker verification (see Sebestyen, 1962). The Sebestyen classifier has a number of properties that makes it useful for the type of pattern recognition problem being considered. These properties will be investigated in later chapters. In implementing the Sebestyen classifier we assume that the ideal reference patterns are known; $M_j=1$, i.e., the ideal reference patterns considered are equally weighted; and $Q_j=(1/2\sigma^2)$.

In the case where only amplitudes are used as features, the univariate noisy measurement can be expressed as

$$x = \left[(A + \nu_1)^2 + \nu_2^2 \right]^{1/2}$$

where A is the uncorrupted value of the amplitude,

$\nu_1, \nu_2 \sim n(0, \sigma^2)$. The probability density function of x , known as the Rician density function, has the functional

form

$$p(x) = (x/\sigma^2) \exp \left[-(1/2\sigma^2)(x^2 + A^2) \right] \cdot I_0(Ax/\sigma^2)$$

where $I_0(z)$, the modified Bessel function of zero order with argument z , is given by the summation

$$I_0(z) = \sum_{n=0}^{\infty} \frac{z^{2n}}{2^{2n} (n!)^2}.$$

The d -variate probability density function of the pattern vector can be written as

$$\begin{aligned} (3-4) \quad p(\underline{x}/\omega_i, \emptyset, \theta) &= \sum_{j=1}^d (x_j/\sigma^2) I_0(f_{ij}(\emptyset, \theta) x_j/\sigma^2) \cdot \\ &\cdot \exp \left[-(1/2\sigma^2) \sum_{j=1}^d (x_j^2 + f_{ij}^2(\emptyset, \theta)) \right]. \end{aligned}$$

The Bayes discriminant function in the amplitude alone case is thus

$$(3-5) \quad g_i(\underline{x}, \phi_o, \theta_o) = P(\omega_i) \sum_{\phi=\phi^Y-a}^{\phi^Y+a} \sum_{\theta=\theta^Z-b}^{\theta^Z+b} \exp \left\{ \sum_{j=1}^d \ln \left[I_o(f_{ij}(\phi, \theta) x_j / \sigma^2) \right] - (1/2 \sigma^2) (x_j^2 + f_{ij}^2(\phi, \theta)) \right\}.$$

When the signal-to-noise ratio is high, the Bayes classifier based on class-conditional densities that are Rician mixtures is closely approximated by the Sebestyen classifier. This can be seen as follows. A high signal-to-noise ratio implies that $f_{ij}(\phi, \theta) \gg \sigma$, which implies that $I_o[f_{ij}(\phi, \theta) x_j / \sigma^2]$ can be approximated by $\exp[f_{ij}(\phi, \theta) x_j / \sigma^2]$. That is, the discriminant function defined by Eq. (3-5) can be approximated by Eq. (3-3).

III.C. Computer Implementation of the Bayes Classifier

The Bayes discriminant function in each of the classification situations considered above are sums of weighted exponentials. If a discriminant function has the form

$$g = \sum_{i=1}^s t_i \exp(a_i)$$

and the exponents a_i can assume an unlimited range of real values, the discriminant functions can easily become either very large or very small and cause overflow or underflow in

any computer implementation. For example, when the Sebestyen classifier is used in the univariate case, the noisy measurement differs from one of the stored values by k units, and the assumed noise standard deviation is σ , then the exponent due to the stored value is $-k^2/2\sigma^2$. If the range of directly computable exponentials on the computer is $\pm E$, then $k^2/2\sigma^2 > E$ implies that g cannot be directly computed. In particular, on the IBM 370/165, $E=174.7$ due to which values of $\sigma=0.2$, $k=4$ are sufficient for this situation to occur.

Since the Bayes classifier only needs to compare discriminant functions and not necessarily compute their absolute values we give here an algorithm that not only eliminates the overflow problem but also avoids computing exponentials in a large number of cases. The algorithm is based on the fact that if the difference between the largest exponentials of each class is greater than a threshold, then the larger sum is the one with the larger maximum exponent. The threshold is determined by the number of terms in each sum. If the difference between the largest exponents does not exceed the threshold, then one of the exponents in each sum is altered and the same test applied again. This procedure is continued until a decision is made.

$$\text{Let } g_1 = \sum_{i=1}^m t_i \exp(a_i) \text{ and } g_2 = \sum_{i=1}^n r_i \exp(b_i)$$

where a_i, b_i are real and $t_i, r_i > 0$. If we assume that the exponent sets $\{a_i\}$ and $\{b_i\}$ are sorted in decreasing order, then the sufficient condition for g_1 being greater than or less than g_2 is as follows:

$$\text{if } a_1 - b_1 \geq \ln(nr_1/t_1) \text{ then } g_1 \geq g_2,$$

$$\text{if } b_1 - a_1 \geq \ln(mt_1/r_1) \text{ then } g_2 \geq g_1$$

where t_1 and r_1 are the weights associated with a_1 and b_1 .

Proof:

Since $g_1 \geq t_1 \exp[a_1]$ and $g_2 \leq n r \exp[b_1]$ it follows that $g_1 \geq g_2 (t_1/nr_1) \exp[a_1 - b_1]$.

Thus, $(t/nr_1) \cdot \exp[a_1 - b_1] \geq 1$ implies that $g_1 \geq g_2$, i.e., $a_1 - b_1 \geq \ln(nr_1/t_1)$ implies that $g_1 \geq g_2$. Similarly we can show $b_1 - a_1 \geq \ln(mt_1/r_1)$ implies that $g_2 \geq g_1$, by considering $g_2 \geq r_1 \cdot \exp[b_1]$ and $g_1 \leq m \cdot t_1 \cdot \exp[a_1]$. Q.E.D.

In a case where the sufficient condition is not satisfied, which happens if

$$(3-6) \ln(r_1/t_1 m) \leq (a_1 - b_1) \leq \ln(nr_1/t_1),$$

we proceed as follows. The two sums of weighted exponentials

can be written as

$$g_1 = t_1 \exp(a_1) + \sum_{i=2}^m t_i \exp(a_i)$$

$$g_2 = r_1 \exp(b_1) + \sum_{i=2}^n r_i \exp(b_i),$$

In order to determine which of the g_i is greater it is enough to determine whether the following inequality holds:

$$t_1 \exp(a_1) - r_1 \exp(b_1) + \sum_{i=2}^m t_i \exp(a_i) > \sum_{i=2}^n r_i \exp(b_i) .$$

Now $[t_1 \exp(a_1) - r_1 \exp(b_1)]$ can be written as $p \cdot \exp(b_1)$ where $p = t_1 \exp(a_1 - b_1) - r_1$. Thus we have to compare two sums of weighted exponentials, one of which has m terms, and the other has $(n-1)$ terms.

As long as r_1 and t_1 are comparable values, the quantity p , which involves $\exp(a_1 - b_1)$, is easily computed. This follows from the bounds on $(a_1 - b_1)$ provided by inequality (3-6). According to the inequality, if $\pm E$ is the range of directly computable exponentials on the computer, then $\exp(a_1 - b_1)$ is directly computable as long as

$$m \cdot \exp(-E) < \frac{r_1}{t_1} < \frac{1}{n} \cdot \exp E. \quad \text{If the weights are not comparable}$$

we can start out by replacing a_i by $\{a_i + \ln t_i\}$, and

b_i by $\{b_i + \ln r_i\}$, which makes all weights equal to unity.

We can thus state the algorithm.

Let $\{a_i\}$, $\{b_i\}$ be the exponent sets associated with g_1 and g_2 respectively. The elements of each of the exponent sets are assumed to be in decreasing order. Let $\{t_i\}$ and $\{r_i\}$ be the corresponding weight sets.

Algorithm B: Bayes discriminant comparison

- C1: Replace each of the a_i by $a_i + \ln t_i$
 Replace each of the b_i by $b_i + \ln r_i$;
- C2: Order each of the sets $\{a_i\}$ and $\{b_i\}$ in
 decreasing order;
- C3: If $a_1 - b_1 \geq \ln n$ choose g_1 and stop;
- C4: If $b_1 - a_1 \geq \ln m$ choose g_2 and stop;
- C5: Replace a_1 by $b_1 + \ln [\exp(a_1 - b_1) - 1]$;
- C6: Delete b_1 from $\{b_i\}$ and set $n = n - 1$;
- C7: Go to C2.

Analysis of the algorithm

The algorithm converges since at the last iteration $n=1$ and the algorithm chooses g_1 . If a decision is not made at either Step C3 or Step C4, the exponential $\exp(a_1 - b_1)$

needs to be computed at Step C5. The maximum and minimum values that $(a_1 - b_1)$ can take at Step C5 are $\ln(n)$ and $\ln(m)$ respectively; if $(a_1 - b_1)$ were outside this range, then a decision would have been made earlier. That is, if the implementation is on a computer whose range of directly computable exponentials is $\pm E$ and $\max(m, n)$ is not greater than $\exp E$, then the overflow/underflow problem is eliminated.

Assuming the weights a_i, b_i to be equal to unity for $i=1, \dots, s$, the computational efficiency of Algorithm B, or a slightly modified version of it, can be compared to the direct procedure of computing $2s$ exponentials and making one comparison. In the best case, a modified form of Algorithm B requires only $2s+1$ comparisons in order to make a decision; the modification is one of determining the maximum elements of sets in Step C2 instead of ordering the sets. In the worst case, s iterations of the algorithm are made, s exponentials are computed at Step C5 and the number of comparisons made is of the order of s .

If Algorithm B is used to implement the Sebestyen classifier, each of the a_i and b_i correspond to a scalar, Q , times the Euclidean distance between the observed pattern \underline{x}

and each of the stored reference patterns corresponding to discrete subclasses. Whenever the Sebestyen classifier makes a decision on the first iteration of the algorithm, it effectively chooses that class which has a reference pattern nearest to \underline{x} . If we assume s discrete subclasses in each class, then the Sebestyen classifier goes through more than one iteration of the algorithm when the difference in the distance from \underline{x} to the nearest pattern lies in the interval $\pm(1/Q)\ln s$. The size of this interval increases with the number of orientations, s . Thus, as s increases the Sebestyen classifier relies less upon the sole evidence of the nearest patterns of each class to \underline{x} .

It is useful to estimate the probability of computing any exponentials when Algorithm B is used to implement the Sebestyen classifier. Consider the univariate problem where noise is Gaussian, i.e., the Bayes classifier is the Sebestyen classifier. Let α and β be equally likely reference patterns of each class such that $\beta > \alpha$. Let the noisy pattern be $x \sim n(\mu, \sigma^2)$ where $\mu = \alpha$ or β . Then the Sebestyen classifier will not make a decision on the first iteration if

$$\ln(1/s) < \frac{(x-\alpha)^2 - (x-\beta)^2}{2\sigma^2} < \ln s \quad ,$$

i.e.,

$$\left[(\alpha + \beta) / 2 \right] - \left[(\sigma^2 \ln s) / (\beta - \alpha) \right] < x < \left[(\alpha + \beta) / 2 \right] + \left[(\sigma^2 \ln s) / (\beta - \alpha) \right]$$

That is, the probability of the Sebestyen classifier computing at least one exponential using Algorithm B is given by

$$\begin{aligned} P_1 &= \text{Prob (compute at least one exponential)} \\ &= \frac{1}{2} \left[P(x/\alpha) + P(x/\beta) \right] \end{aligned}$$

where

$$P(x/\mu) = \int_{\frac{\alpha + \beta}{2} - K}^{\frac{\alpha + \beta}{2} + K} p(x/\mu) \cdot dx$$

$$p(x/\mu) = n(\mu, \sigma^2), \text{ and } K = \frac{\sigma^2 \ln s}{(\beta - \alpha)}.$$

The above expressions can be used to obtain conditions on the relative locations of α and β for which P_1 is less than a specified value. For example, P_1 is less than 0.05, which happens when

$$\frac{\alpha + \beta}{2} - K \geq \alpha + 2\sigma, \text{ and } \frac{\alpha + \beta}{2} + K \leq \beta - 2\sigma.$$

That is, P_1 is less than 0.05 if $\frac{\beta - \alpha}{2} \geq 2\sigma + K$. It follows

that for values of s that are less than $\exp(8)$, P_1 is less than 0.05 provided $(\beta - \alpha) > 4\sigma$.

The above analysis has made the simplifying assumption that α and β are the closest reference points to x , and that the problem is univariate. However it indicates the fact that the probability of computing exponentials using Algorithm B is directly dependent upon the noise level, and inversely dependent upon the interclass separation. In the multivariate situation the probability of computing exponentials was determined experimentally with data from the radar aircraft identification problem. It was found that when computer simulation experiments were conducted with $d=8$, $s=9$, and $\sigma=0.5$, in 60% of the experiments a decision was made on the first iteration of the algorithm.

III.D. The Bayes Misclassification Probability

The most effective measure of the performance of a classification algorithm is the probability of misclassification associated with the algorithm. This probability is generally taken to be the average misclassification probability, where the average is obtained by considering all possible classes and all possible values of the pattern vector. In this section we examine the effect of considering

multiple orientation possibilities on the Bayes misclassification probability, P^* . An approximation to P^* will then be obtained, and its dependency on the location of reference patterns will be examined.

Let P^* be the misclassification probability of the Bayes classifier when s discrete orientations are equally likely. At each of the s orientations let the Bayes misclassification probability be P_j^* , $j=1, \dots, s$. It will be shown that in general P^* is greater than the average value of P_j^* , where the average is over s values.

The Bayes misclassification probability between two classes can be expressed as (see Chen, 1972)

$$(3-7) \quad P^* = \frac{1}{2} K(f_1, f_2)$$

where $K(f_1, f_2)$, the Kolmogorov variational distance between classes with class-conditional probability densities f_1, f_2 , is given by

$$(3-8) \quad K(f_1, f_2) = \frac{1}{2} \int |P(\omega_1) f_1 - P(\omega_2) f_2| \cdot d\mathbf{x} \quad ,$$

$P(\omega_1)$, $P(\omega_2)$ being the a priori probabilities of classes.

If f_{ij} denotes the probability density function associated with the j^{th} orientation of the i^{th} class, then the

Bayes misclassification probability between the two classes at the j^{th} orientation is given by

$$p_j^* = \frac{1}{2} - K(f_{1j}, f_{2j}).$$

Since the probability density function associated with class ω_i when s discrete orientation are equally likely is given by

$$f_i = (1/s) \sum_{j=1}^s f_{ij},$$

it follows that the Bayes misclassification probability of a classifier that considers s orientations to be equally likely has the form

$$(3-9) \quad P^* = \frac{1}{2} - K \left[(1/s) \sum_{j=1}^s f_{1j}, (1/s) \sum_{j=1}^s f_{2j} \right].$$

Since the Kolmogorov variational distance between f_{1j} and f_{2j} has the form

$$K(f_{1j}, f_{2j}) = \frac{1}{2} \int |P(\omega_1) f_{1j} - P(\omega_2) f_{2j}| \cdot d\mathbf{x},$$

and the fact that for real a_j, b_j we can write

$$\left| \sum_{j=1}^s a_j - \sum_{j=1}^s b_j \right| \leq \sum_{j=1}^s |a_j - b_j|,$$

it follows that

$$(3-10) \quad K \left[(1/s) \sum_{j=1}^s f_{1j}, (1/s) \sum_{j=1}^s f_{2j} \right] \leq (1/s) \sum_{j=1}^s K(f_{1j}, f_{2j}).$$

Combining (3-9) and (3-10) we can write

$$P^* \geq \frac{1}{2} - (1/s) \sum_{j=1}^s K(f_{1j}, f_{2j}),$$

or,

$$P^* \geq (1/s) \sum_{j=1}^s \left[\frac{1}{2} - K(f_{1j}, f_{2j}) \right],$$

or

$$(3-11) \quad P^* \geq (1/s) \sum_{j=1}^s p_j^*.$$

A different proof of inequality (3-11) for the case of $s=2$ is given by Kanal and Chandrasekaran (1972). They obtain conditions for the strict inequality to hold for the univariate, Gaussian case.

For arbitrary class-conditional probability density functions, including the ones encountered in the radar aircraft identification problem, it is difficult to compute P^* by evaluating Eq. (3-7). However, bounds on P^* for exponential distributions can be computed easily in terms of the Bhattacharyya coefficient. The Bhattacharyya coefficient between classes ω_1 and ω_2 is defined by

$$(3-12) \quad \rho = \int \left[p(\underline{x}/\omega_1) \cdot p(\underline{x}/\omega_2) \right]^{\frac{1}{2}} \cdot d\underline{x}.$$

If $p(\underline{x}/\omega_1)$ and $p(\underline{x}/\omega_2)$ are exponential functions then their product is an exponential function which makes the integral of Eq. (3-12) simpler to evaluate than the integral of Eq. (3-8).

The bounds on P^* in terms of ρ have been shown to be

$$(3-13) \quad \frac{1}{2}\rho \geq [P(\omega_1)P(\omega_2)]^{\frac{1}{2}} \rho \geq P^* \geq \frac{1}{2}[1-(1-\rho^2)^{\frac{1}{2}}] \geq \frac{1}{4}\rho^2$$

Concerning the tightness of the bounds, it is known that the upper bound $\frac{1}{2}\rho$ is "good, though no general statement is possible" (Kailath, 1967).

If each of the class-conditional densities are Gaussian mixtures, which is the case in the radar aircraft identification problem when both amplitude and phase are used to construct the pattern vector, then $p(\underline{x}/\omega_i)$ has the form

$$p(\underline{x}/\omega_i) = (1/s) \sum_{j=1}^s \frac{1}{(2\pi)^{d/2} \sigma^d} \exp \left\{ - \|\underline{x} - \underline{a}_{ij}\|^2 / 2\sigma^2 \right\}$$

where \underline{a}_{ij} is the ideal d -dimensional pattern vector for the j^{th} orientation of the i^{th} class, $\|\underline{x} - \underline{a}_{ij}\|$ is the Euclidean norm of vector $\underline{x} - \underline{a}_{ij}$, and s discrete orientations are considered to be equally likely. In this case, the Bhattacharyya coefficient between the two classes can be

written as

$$\rho = \frac{1}{(2\pi)^{d/2} \sigma^d s^2} \sum_{j=1}^s \sum_{k=1}^s \int_{-\infty}^{\infty} \exp \left\{ - \left[\|\underline{x} - \underline{a}_{1j}\|^2 + \|\underline{x} - \underline{a}_{2k}\|^2 \right] / 4 \sigma^2 \right\} d\underline{x} .$$

After some algebra, the exponent can be written as

$$\left\{ - \left\| \underline{x} - \frac{\underline{a}_{1j} + \underline{a}_{2k}}{2} \right\|^2 / 2 \sigma^2 \right\} - \left\{ \left\| \underline{a}_{1j} - \underline{a}_{2k} \right\|^2 / 8 \sigma^2 \right\} .$$

By observing that the first term is the exponent of the probability density function, $n \left[(\underline{a}_{1j} + \underline{a}_{2k}) / 2, \sigma^2 I \right]$, the Bhattacharyya coefficient reduces to

$$(3-14) \quad \rho = (1/s^2) \sum_{j=1}^s \sum_{k=1}^s \exp \left\{ - \left\| \underline{a}_{1j} - \underline{a}_{2k} \right\|^2 / 8 \sigma^2 \right\} .$$

Thus the Bhattacharyya coefficient between two classes over a specified set of orientations can be obtained by computing the Euclidean distances between patterns of each class, and by obtaining the average value of the exponential of each of the distances weighted by $(-1/8\sigma^2)$.

The Bhattacharyya coefficient has certain properties that make it useful for a classification problem where classes are characterized by ideal reference patterns, and noise is Gaussian. It can be used to obtain bounds on the Bayes misclassification probability in terms of the maximum and minimum inter-class separation of patterns. It can also

be used to predict the behaviour of the Bayes misclassification probability with changes in the value of noise standard deviation. Each of these aspects of the Bhattacharyya coefficient will be examined in the rest of this section.

Let $\{\underline{a}_{1j}\}$ and $\{\underline{a}_{2k}\}$ be two sets of reference patterns that correspond to classes ω_1 and ω_2 respectively. We will define λ_1 and λ_2 to be the minimum and maximum interclass separation of reference patterns. In other words,

$$\lambda_1 = \min_{j,k} \left\{ \|\underline{a}_{1j} - \underline{a}_{2k}\|^2 \right\}, \quad \lambda_2 = \max_{j,k} \left\{ \|\underline{a}_{1j} - \underline{a}_{2k}\|^2 \right\}.$$

Let ρ_{jk} be the Bhattacharyya coefficient between subclasses associated with \underline{a}_{1j} and \underline{a}_{2k} , i.e.,

$$\rho_{jk} = \exp \left\{ -\|\underline{a}_{1j} - \underline{a}_{2k}\|^2 / 8\sigma^2 \right\}.$$

In terms of these definitions we can obtain a bound on the Bayes misclassification probability in terms of maximum and minimum inter-class separations as follows. Since

$$\rho = (1/s^2) \sum_{j=1}^s \sum_{k=1}^s \rho_{jk}, \quad \max_{j,k} \rho_{jk} = \exp \left[-\lambda_1 / 8\sigma^2 \right],$$

$$\text{and } \min_{j,k} \rho_{jk} = \exp \left[-\lambda_2 / 8\sigma^2 \right],$$

it follows that

$$(3-15) \quad \exp \left[-\lambda_2 / 8\sigma^2 \right] \leq \rho \leq \exp \left[-\lambda_1 / 8\sigma^2 \right].$$

Combining inequalities (3-13) and (3-15) we get

$$(3-16) \quad \frac{1}{2} \cdot \exp(-\lambda_1/8\sigma^2) \geq P^* \geq \frac{1}{4} \cdot \exp(-\lambda_2/4\sigma^2) \quad .$$

Inequality (3-16) provides a simple bound on the Bayes misclassification probability when classes are characterized by ideal reference patterns, and noise is Gaussian. In order to compute these bounds, only the noise variance, and minimum and maximum inter-class pattern separation need to be known.

The bounds on P^* , as provided by (3-16) tend to get loose as the difference between λ_1 and λ_2 is increased. For specified values of λ_1 and λ_2 , the bound provided by (3-16) is tightest when σ^2 has the value σ_t^2 defined by

$$(3-17) \quad \sigma_t^2 = \frac{2 \cdot \lambda_2 - \lambda_1}{8 \ln(\lambda_2/\lambda_1)} \quad .$$

Condition (3-17) was obtained by determining the value of σ^2 which minimizes the difference between the upper and lower bounds on P^* in (3-16).

In general, the bounds provided by inequality (3-16) are loose. If reference patterns are widely scattered, which makes λ_1 and λ_2 very different, then P^* can be approximated by the upper bound provided by ρ . That is,

$$(3-18) \quad P^* \cong \frac{1}{2} \rho \quad .$$

In order to determine the goodness of approximation (3-18),

computed values of $\frac{1}{2}\rho$ were compared to estimates of P^* obtained by Monte Carlo simulation. With $d=8$, $s=9$, and data from the radar aircraft identification problem, it was found that the relative error in the approximation does not exceed 4%, thereby demonstrating the utility of the approximation to the problem.

In implementing a classifier for problems such as radar aircraft identification, the performance of the classifier with variation in the noise level is of interest. We examine here the behavior of ρ , as defined by Eq. (3-14) with changes in the noise standard deviation σ . It will be shown that ρ is more sensitive to changes in σ when $\sigma < \sigma_1$ than when $\sigma > \sigma_2$, where σ_1 and σ_2 depend upon the configuration of reference patterns in feature-space. The actual values of σ_1 and σ_2 will be shown lie in the ranges defined by the following inequality

$$(3-19) \quad (2/3) \cdot \lambda_1 \leq \sigma_1^2 \leq \sigma_2^2 \leq (2/3) \cdot \lambda_2$$

If the ratio (λ_2/λ_1) is less than the constant 1.86, then the values of σ_1 and σ_2 will be shown to be identical. This implies that when the maximum to minimum inter-class separation is less than 1.86, then the ρ versus σ plot has

a well-defined knee. In other words, the rate of increase of ρ with σ is unimodal, with a peak at $\sigma_0 = \sigma_1 = \sigma_2$. This behavior of ρ with σ is illustrated in Fig. 3-1. The Bhattacharyya coefficient between two classes is a linear combination of the Bhattacharyya coefficient between pairs of reference patterns with unlike classification. If ρ_{jk} is the Bhattacharyya coefficient between subclasses associated with a_{1j} and a_{2k} , then $d\rho_{jk}/d\sigma$ is a unimodal function of σ . It will be seen that if $(\lambda_2/\lambda_1) < 1.86$ then each of the $d\rho_{jk}/d\sigma$ have peaks that are close enough to make the overall derivative, $d\rho/d\sigma$, a unimodal function of σ .

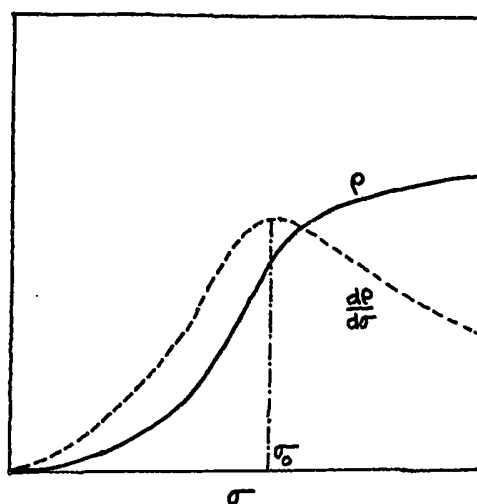


Fig. 3-1. Behavior of Bhattacharyya coefficient with σ when $(\lambda_2/\lambda_1) < 1.86$.

Theorem 3.1 Let a_{ij} be the reference patterns corresponding to class ω_i at the j^{th} orientation, where $i=1, 2, j=1, \dots, s$.

Let the noisy pattern be distributed according to $n(\underline{a}_{ij}, \sigma^2)$. Let λ_1 and λ_2 be the minimum and maximum inter-class separation of reference patterns. If each of the s orientations is equally likely, then the Bhattacharyya coefficient between the two classes has the following properties:

- (i) ρ is monotone increasing with σ ;
- (ii) the rate of increase of ρ with σ is monotone increasing for $\sigma < \sigma_1$, is monotone decreasing for $\sigma > \sigma_2$, where

$$\sigma_1, \sigma_2 \in \left[(2/3) \cdot \lambda_1, (2/3) \cdot \lambda_2 \right];$$
- (iii) if $(\lambda_2 / \lambda_1) < 1.86$ then $\sigma_1 = \sigma_2$.

Proof: Let ω_1 and ω_2 be the two classes, and let $A_{jk} = \|\underline{a}_{1j} - \underline{a}_{2k}\|^2$. The Bhattacharyya coefficient between ω_1 and ω_2 is defined by Eq. (3-14). The first derivative of ρ with respect to σ is

$$\frac{d\rho}{d\sigma} = \left[\frac{2}{s^2 \sigma^3} \right] \cdot \sum_{j=1}^s \sum_{k=1}^s A_{jk} \cdot \exp(-A_{jk}/8\sigma^2)$$

which is positive for all positive values of σ , and so ρ is monotone increasing with σ . The second derivative of ρ with respect to σ is

$$\frac{d^2\rho}{d\sigma^2} = \left[\frac{2}{s^2 \sigma^4} \right] \cdot \sum_{j=1}^s \sum_{k=1}^s A_{jk} \cdot \left[(2A_{jk}/\sigma^2) - 3 \right] \cdot \exp(-A_{jk}/8\sigma^2)$$

If $\sigma^2 < (2/3) \min_{j,k} \{A_{jk}\} = (2/3) \cdot \lambda_1$, then $d^2\rho/d\sigma^2$ is positive since every one of its component terms will be positive. Similarly, if $\sigma^2 > (2/3) \max_{j,k} \{A_{jk}\} = (2/3) \cdot \lambda_2$, then $d^2\rho/d\sigma^2$ is negative. That is, there is a value of σ_1^2 , which is at least $(2/3) \cdot \lambda_1$, such that $d\rho/d\sigma$ is monotone increasing for $\sigma^2 < \sigma_1^2$. Similarly, there is a value of σ_2^2 which is at most $(2/3) \cdot \lambda_2$ such that $d\rho/d\sigma$ is monotone decreasing for $\sigma^2 > \sigma_2^2$.

We will show that $\sigma_1 = \sigma_2$ for $(\lambda_2/\lambda_1) < 1.86$. Consider the function $f(\sigma^2)$ which has two additive components, each component having the form of the components of ρ as defined by Eq. (3-14). That is,

$$f(\sigma^2) = f_1(\sigma^2) + f_2(\sigma^2)$$

where $f_i(\sigma^2) = \exp(-A_i/8\sigma^2)$. The derivatives of $f_i(\sigma^2)$ with respect to σ are

$$\begin{aligned} f_i'(\sigma^2) &= (2/\sigma^3) \cdot A_i \cdot \exp(-A_i/8\sigma^2) \\ f_i''(\sigma^2) &= (2/\sigma^4) \cdot A_i \cdot [(2A_i/\sigma^2) - 3] \cdot \exp(-A_i/8\sigma^2) \\ f_i'''(\sigma^2) &= (4/\sigma^5) \cdot A_i \cdot [(2A_i^2/\sigma^4) - (9A_i/\sigma^2) + 6] \\ &\quad \cdot \exp(-A_i/8\sigma^2). \end{aligned}$$

The expression $[(2A_i^2/\sigma^4) - (9A_i/\sigma^2) + 6]$ is equal to zero for $\sigma^2 = 0.27A_i$ and for $\sigma^2 = 1.23A_i$; the values were obtained by setting the expression equal to zero and

solving for σ^2 . The expression, and hence $f_1'''(\sigma^2)$, is positive outside the interval $(0.27A_1, 1.23A_1)$ and is negative within the interval. Since $f_1''(\sigma^2)$ is positive when $\sigma^2 < 0.66A_1$ and negative otherwise, the behavior of $f_1''(\sigma^2)$ with σ^2 is as follows: it starts with a positive value, keeps increasing until a maximum is reached at $0.27A_1$, after which it starts to decrease, is zero at $0.66A_1$, reaches a minimum at $1.23A_1$, starts to increase, and asymptotically reaches zero. From the expression for $f_1'(\sigma^2)$ and the above behavior of $f_1''(\sigma^2)$ it follows that $f_1'(\sigma^2)$ is always positive, and is unimodal with a peak at $0.66A_1$. The behavior of $f_1''(\sigma^2)$ and $f_1'(\sigma^2)$ with σ^2 are shown in Fig. 3-2.

We will consider now the behavior of the function $f''(\sigma^2) = f_1''(\sigma^2) + f_2''(\sigma^2)$ and show that if $1 < (A_2/A_1) < 1.86$, then $f''(\sigma^2)$ is unimodal. $(A_2/A_1) > 1$ implies that the zeros of $f''(\sigma^2)$ can only be in the interval $(0.66A_1, 0.66A_2)$, due to the fact that $f_1''(\sigma^2)$ and $f_2''(\sigma^2)$ have the same sign outside the interval. Also there exists at least one zero of $f''(\sigma^2)$ since it is positive for $\sigma^2 > 0.66A_1$ and negative for $\sigma^2 < 0.66A_2$. Now, $A_2/A_1 < 1.86$ implies that $0.66A_2 < 1.23A_1$ and $0.27A_2 < 0.66A_1$; that is $f_1''(\sigma^2)$ and $f_2''(\sigma^2)$ are both

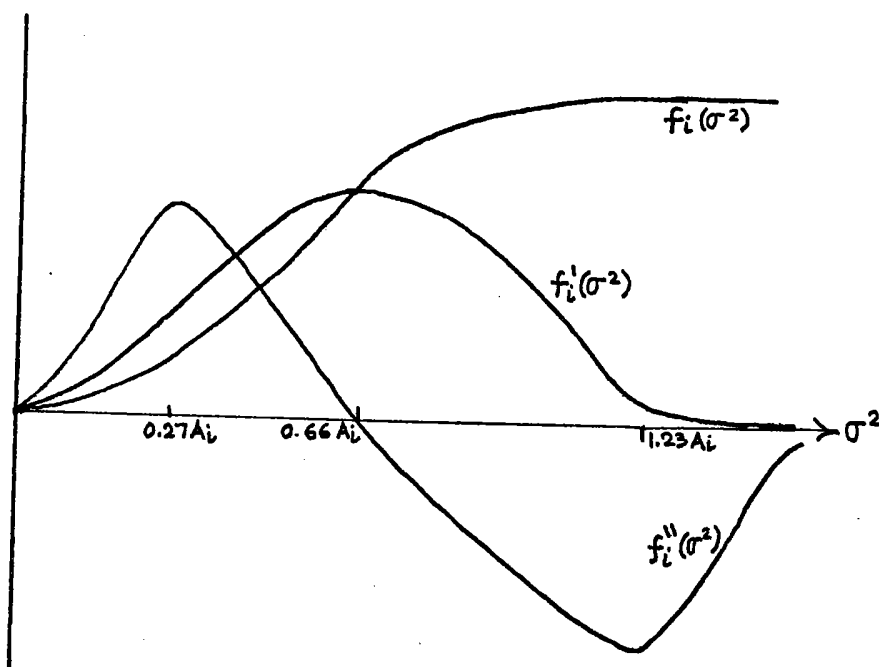


Fig. 3-2. Behavior of derivatives of $f_i(\sigma^2)$ with σ^2 .

monotone decreasing in $(0.66A_1, 0.66A_2)$. It follows that there exists a unique zero of $f''(\sigma^2)$ in $(0.66A_1, 0.66A_2)$.

In other words, $f'(\sigma^2)$ is unimodal with a peak at

$\sigma_0^2 \in (0.66A_1, 0.66A_2)$. Thus the rate of increase of $f(\sigma^2)$ is monotone increasing for $\sigma^2 < \sigma_0^2$ and monotone decreasing for $\sigma^2 > \sigma_0^2$.

By adding a third component $f_3(\sigma^2)$ to $f(\sigma^2)$ we can show that $f'(\sigma^2)$ is unimodal if $1 < (A_{\max}/A_{\min}) < 1.86$, as follows. Let $f_3(\sigma^2) = \exp(-A_3/8\sigma^2)$, and without loss of

generality let $A_1 < A_3 < A_2$. Since $(A_2/A_1) < 1.86$, it follows that $0.27A_3 < 0.66A_1$, $0.66A_2 < 1.23A_3$. Thus $f_3''(\sigma^2)$ is monotone decreasing in $(0.66A_1, 0.66A_2)$ and has a zero in the interval. As a consequence $f''(\sigma^2) = f_1''(\sigma^2) + f_2''(\sigma^2) + f_3''(\sigma^2)$ has a unique zero in $(0.66A_1, 0.66A_2)$. The result is likewise extended to the case where $f(\sigma^2)$ has arbitrary number of additive components.

The Bhattacharyya coefficient, ρ , can be verified from Eq. (3-14) to have the form $(1/s^2) \cdot f(\sigma^2)$, where $f(\sigma^2)$ has s^2 additive components. Using the above argument, $d^2\rho/d\sigma^2$ has a unique zero if $(\lambda_2/\lambda_1) < 1.86$. It follows that as long as $(\lambda_2/\lambda_1) < 1.86$, there is a unique value of σ_0 such that $\sigma_1 = \sigma_2 = \sigma_0$.

Q.E.D.

The relevance of constant 1.86 to the above analysis can be restated as follows. The component of $d^2\rho/d\sigma^2$ due to λ_1 has a maximum at $0.27 \cdot \lambda_1$, is zero at $0.66 \cdot \lambda_1$, and has a minimum at $1.23 \cdot \lambda_1$. The corresponding values of the component due to λ_2 are $0.27 \cdot \lambda_2$, $0.66 \cdot \lambda_2$, and $1.23 \cdot \lambda_2$ respectively. If $0.66 \cdot \lambda_2 < 1.23 \cdot \lambda_1$ and $0.27 \cdot \lambda_2 < 0.66 \cdot \lambda_1$, then $d^2\rho/d\sigma^2$ has a unique zero. In other words, $(\lambda_2/\lambda_1) < 1.86$

implies that there is a unique value of σ at which the second derivative of ρ with respect to σ is equal to zero.

The values of ρ at $\sigma^2 = (2/3)\lambda_1$ and $(2/3)\lambda_2$ can be bounded using inequality (3-15). Writing ρ as a function of σ , or as $\rho(\sigma)$, and using (3-15) we have

$$(3-20) \quad \begin{aligned} \exp(-3\lambda_2/16\lambda_1) &\leq \rho[\sqrt{(2/3)\lambda_1}] \leq 0.84 \\ 0.84 &\leq \rho[\sqrt{(2/3)\lambda_2}] \leq \exp(-3\lambda_1/16\lambda_2) \end{aligned}$$

In the particular case where $(\lambda_2/\lambda_1) \leq 1.86$, inequalities (3-20) reduce to

$$(3-21) \quad \begin{aligned} 0.71 &\leq \rho[\sqrt{(2/3)\lambda_1}] \leq 0.84 \\ 0.84 &\leq \rho[\sqrt{(2/3)\lambda_2}] \leq 0.90 \end{aligned}$$

or the knee of the ρ versus σ plot which occurs at σ_0 has a value that lies between 0.71 and 0.90. The value of σ_0 at which $d\rho/d\sigma$ is maximum can be determined by the procedure given in Appendix A.

If the ratio of maximum to minimum inter-class separation exceeds 1.86, i.e., $(\lambda_2/\lambda_1) > 1.86$, then there can be more than one zero of $d^2\rho/d\sigma^2$, all of which lie in the interval $(0.66\lambda_1, 0.66\lambda_2)$. Multiple zeros of the second derivative of ρ imply that $d\rho/d\sigma$ is not unimodal. The possibility of multiple zeros of the second derivative can be seen by considering different ranges of values for

(A_2/A_1) , and substituting in the equation $f''(\sigma^2) = f_1''(\sigma^2) + f_2''(\sigma^2)$. The reader is referred to Appendix A for a characterization of where the zeros lie as the ratio (A_2/A_1) is increased.

It follows from the above that anomalous behavior of ρ with σ can be expected in the interval $(0.66\lambda_2, 0.66\lambda_1)$; the anomaly is with respect to the derivatives, ρ as such is known to be monotone increasing with σ . Fig. 3-3 illustrates possible behavior of $\frac{1}{2}\rho$ with σ for the cases where $(\lambda_2/\lambda_1) > 1.86$, and $(\lambda_2/\lambda_1) < 1.86$. We have chosen to consider $\frac{1}{2}\rho$ since it is an approximation to the Bayes misclassification probability.

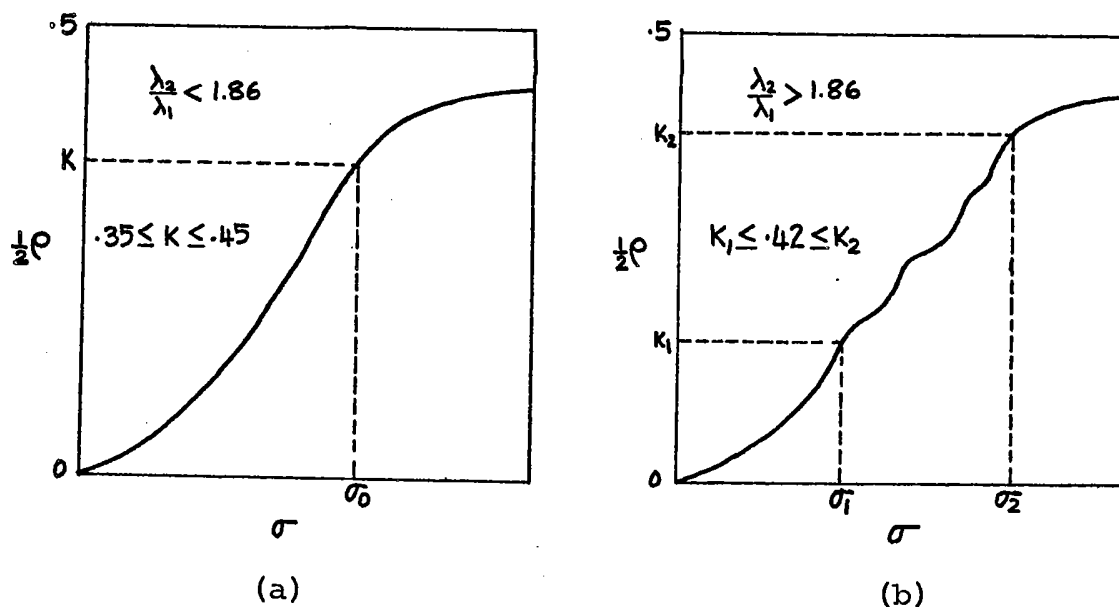


Fig. 3-3 Behavior of $\frac{1}{2}\rho$ as a function of σ .

Inequalities (3-20) and (3-21) provide some insight into the values of $\frac{1}{2}\rho$ in the interval (σ_1, σ_2) . From inequality (3-21) it follows that when $(\lambda_2/\lambda_1) < 1.86$, and $\sigma = \sigma_0$, $\frac{1}{2}\rho$ has a value that lies in the interval $(0.35, 0.45)$. In the case where $(\lambda_2/\lambda_1) > 1.86$, from inequality (3-20) it follows that $\frac{1}{2}\rho$ has a value 0.42 in the interval $\left\{ \left[(2/3) \cdot \lambda_1 \right]^{\frac{1}{2}}, \left[(2/3) \cdot \lambda_2 \right]^{\frac{1}{2}} \right\}$.

The implications of the above analysis to the implementation of a radar aircraft identification system are as follows. The primary use is in the ability to predict analytically the behaviour of average misclassification probability with noise standard deviation. This is based on the fact that the average misclassification probability can be approximated by half the value of the Bhattacharyya coefficient between classes.

Another use of the above analysis is in determining noise levels of importance in evaluating gathered radar aircraft data. If the maximum inter-class separation between classes is small, then the misclassification probability between classes reaches a saturation point at a small value of noise standard deviation. That is, if we are interested

in evaluating the Bayes misclassification probability for a specified set of reference patterns, and at different noise levels, the maximum value of noise standard deviation that needs to be considered is $[(2/3)\lambda_2]^{1/2}$. This follows from the fact that beyond $[(2/3)\lambda_2]^{1/2}$, the rate of increase of misclassification probability asymptotically reaches zero at a rate that is progressively slower. However, if the minimum inter-class separation between classes is large, any small change in the noise level will lead to considerable change in the resulting misclassification probability. If the ratio of the maximum to minimum inter-class separation is less than 1.86, there is a unique knee in the misclassification probability versus noise level plot at the value defined by σ_0 . Thus a noise level of importance in evaluating the gathered data is $\sigma = \sigma_0$. If the ratio exceeds 1.86, then the interesting region of noise-level is between $[(2/3)\lambda_2]^{1/2}$ and $[(2/3)\lambda_1]^{1/2}$ where the anomalous behavior illustrated in Fig. 3-3b can be expected.

CHAPTER IV

THE DECISION SURFACES OF THE SEBESTYEN AND NEAREST-NEIGHBOR CLASSIFIERS

In the previous chapter the classification problem was stated in decision-theoretic terms. For an important case of the radar aircraft identification problem, the Bayes classifier was shown to have the form of the Sebestyen classifier. The Bayes classifier was examined from the following two points of view: efficient implementation on a digital computer, and analytical determination of misclassification probability. In this chapter we will look at the relationship between the Bayes classifier and the family of nearest-neighbor classifiers, when classes are characterized by a deterministic set of reference patterns. We will then examine the dependency of the decision surfaces of the Sebestyen and nearest-neighbor classifiers on:

- (i) assumed value of noise variance, and
- (ii) change in the relative location of reference patterns due to change in dimensionality.

IV.A. Nearest-Neighbor Classification

If every class is represented by a set of sample patterns in feature-space, then the nearest-neighbor rule of classification is as follows. In order to classify a pattern \underline{x} , the distance from \underline{x} to every sample pattern is measured. The classification of \underline{x} is taken to be the classification of the pattern that is nearest, in terms of the measured distance, to \underline{x} .

A classifier based on the nearest-neighbor rule is the nearest-neighbor classifier (NNC) and the set of sample patterns constitute its reference set. The distance measure of the nearest-neighbor classifier can be any one of many distance measures such as the Euclidean distance, Euclidean distance weighted by the a priori probability associated with reference patterns, statistical distance measures based on functional forms of distributions, etc.

The nearest-neighbor classifier with a Euclidean distance measure has been applied to the radar aircraft identification problem (Lin, 1974). In this application, the ideal patterns from all possible aircraft at all possible orientations constitute the reference set of the

nearest-neighbor classifier. Essentially, the pattern \underline{x} is classified into that class whose data surface is nearest to \underline{x} . Since the sample patterns constitute a deterministic set rather than a statistical one, the nearest-neighbor classifier with a Euclidean distance measure is also referred to as the minimum distance classifier here. The standard definition of the minimum distance classifier corresponds to one that assigns \underline{x} to that class ω_i whose ideal pattern, usually the mean of a statistical distribution, is closest in the Euclidean sense to \underline{x} . The performance of the nearest-neighbor classifier has been found to be far better than that of a classifier which attempts to separate every pair of classes by a single hyperplane, i.e., the linear discriminant function classifier. The reason for the poor performance of the linear discriminant classifier is the highly convoluted nature of data surfaces in the radar aircraft identification problem, due to which separating hyperplanes cannot be found. We restrict our comparison of the Sebestyen classifier to the nearest-neighbor classifier.

IV.B. Relationship Between the Bayes and Nearest-Neighbor Classifiers

In order to relate the Bayes classifier of Chapter III to the family of nearest-neighbor classifiers, the classification problem will be viewed here as a problem with classes and subclasses. Every class may be thought of as representing an object, and every subclass as representing a discrete orientation of the object.

Consider the classification problem with c classes $\omega_1, \dots, \omega_c$. Let class ω_i have s subclasses labelled $\omega_{i1}, \dots, \omega_{is}$. The information available to the classifier is in terms of the following probabilities and probability density functions: $P(\omega_i)$, the a priori probability of class ω_i ; $P(\omega_{ij})$, the a priori probability of subclass ω_{ij} given that the true class is ω_i ; $p(\underline{x}/\omega_{ij})$, the probability density function of \underline{x} conditioned on the subclass.

Consider two parametric classifiers C_{mult} and C_{mix} which are designed for the above problem as follows:

C_{mult} treats the problem as a multiclass problem with $c \cdot s$ subclasses, where each subclass is considered to have an a priori probability $P(\omega_{ij})P(\omega_i)$. Utilizing the class-conditional densities $p(\underline{x}/\omega_{ij})$, C_{mult} computes

discriminant functions for each of the c subclasses and classifies \underline{x} into that subclass, and also into that class associated with that subclass, for which the discriminant function is a maximum. That is, C_{mult} chooses that class ω_i for which

$$g_{ij}(\underline{x}) = p(\underline{x}/\omega_{ij}) \cdot P(\omega_{ij}) \cdot P(\omega_i)$$

is maximum over all i and j .

C_{mix} is another parametric classifier for the same problem which computes only c discriminant functions. It is based on the fact that the class-conditional probability density functions can be written as

$$p(\underline{x}/\omega_i) = \sum_{j=1}^S p(\underline{x}/\omega_{ij}) \cdot P(\omega_{ij})$$

Thus C_{mix} computes c discriminant functions of the form

$$g_i(\underline{x}) = \sum_{j=1}^S p(\underline{x}/\omega_{ij}) \cdot P(\omega_{ij}) \cdot P(\omega_i)$$

and chooses that class for which $g_i(\underline{x})$ is maximum. The classifiers C_{mult} and C_{mix} are schematically illustrated in Fig. 4-1.

Since the goal is to choose the class and not the subclass, the Bayes classifier for the problem is one which is based on the overall class-conditional probability density functions $p(\underline{x}/\omega_i)$, i.e., C_{mix} . Since both C_{mix}

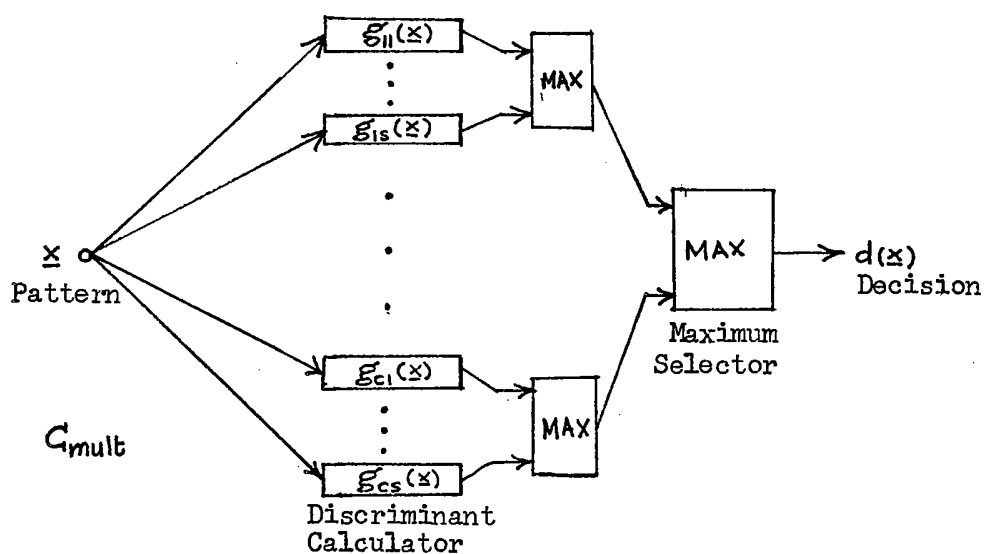
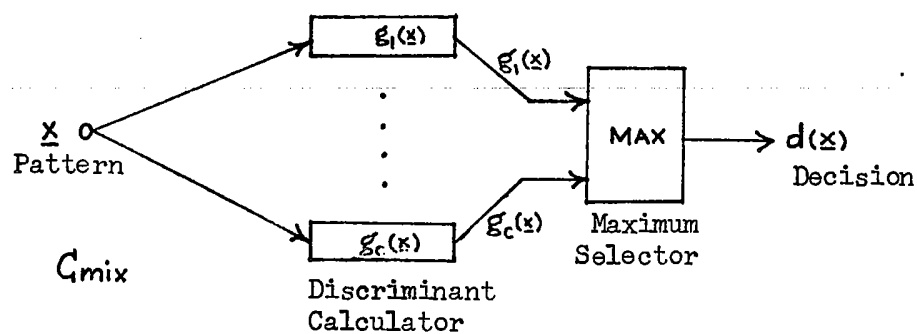


Fig. 4-1 Schematic illustration of classifiers C_{mult} and C_{mix}

and C_{mult} make use of all available information, it might appear that the two classifiers are identical. However, the two classifiers are not identical as can be seen by considering values of $g_{ij}(\underline{x})$ such that

$$\max_j g_{ij}(\underline{x}) > \max_j g_{kj}(\underline{x}) \text{ and } \sum_{j=1}^s g_{ij}(\underline{x}) < \sum_{j=1}^s g_{kj}(\underline{x}) .$$

It is sufficient to show the non-identity of C_{mix} and C_{mult} in a simple case. We can draw the necessary conclusion from this.

Let $p(x/\omega_{ij})$ be univariate Gaussian with mean a_{ij} and standard deviation σ . Let $c=2$, $s=2$, and all $P(\omega_{ij})$ equal. Let the relative locations of a_{11} , $a_{12} \in \omega_1$, and a_{21} , $a_{22} \in \omega_2$ be as shown in Fig. 4-2. Let $a_{22}-a_{12}=a_{12}-a_{21}=\sigma$, and let $a_{21}-a_{11} \gg 3\sigma$. Due to the Gaussian probability density function being symmetric about its mean, we can state that the subclass discriminant functions satisfy the following. Over all subclass discriminant functions $g_{ij}(x)$, $g_{11}(x)$ is maximum for $x \in (-\infty, c_1)$; $g_{21}(x)$ is maximum for $x \in (c_1, c_2)$; $g_{12}(x)$ is maximum for $x \in (c_2, c_3)$; $g_{22}(x)$ is maximum for $x \in (c_3, \infty)$, where $c_1 = \frac{1}{2}(a_{11}+a_{21})$, $c_2 = \frac{1}{2}(a_{21}+a_{12})$, and $c_3 = \frac{1}{2}(a_{12}+a_{22})$. The discriminant functions of C_{mix} are $g_1(x) = \frac{1}{2}[g_{11}(x)+g_{12}(x)]$ and $g_2(x) = \frac{1}{2}[g_{21}(x)+g_{22}(x)]$. The two discriminant functions are illustrated in Fig. 4-2b. Consider now the situation where $x \in (c_2, c_3)$. Since $g_{12}(x)$ is maximum in this interval, C_{mult} chooses class ω_1 . On the other hand $g_2(x) > g_1(x)$ in this interval, thus C_{mix} chooses class ω_2 . That is, classifiers C_{mult} and C_{mix} do not always make identical decisions.

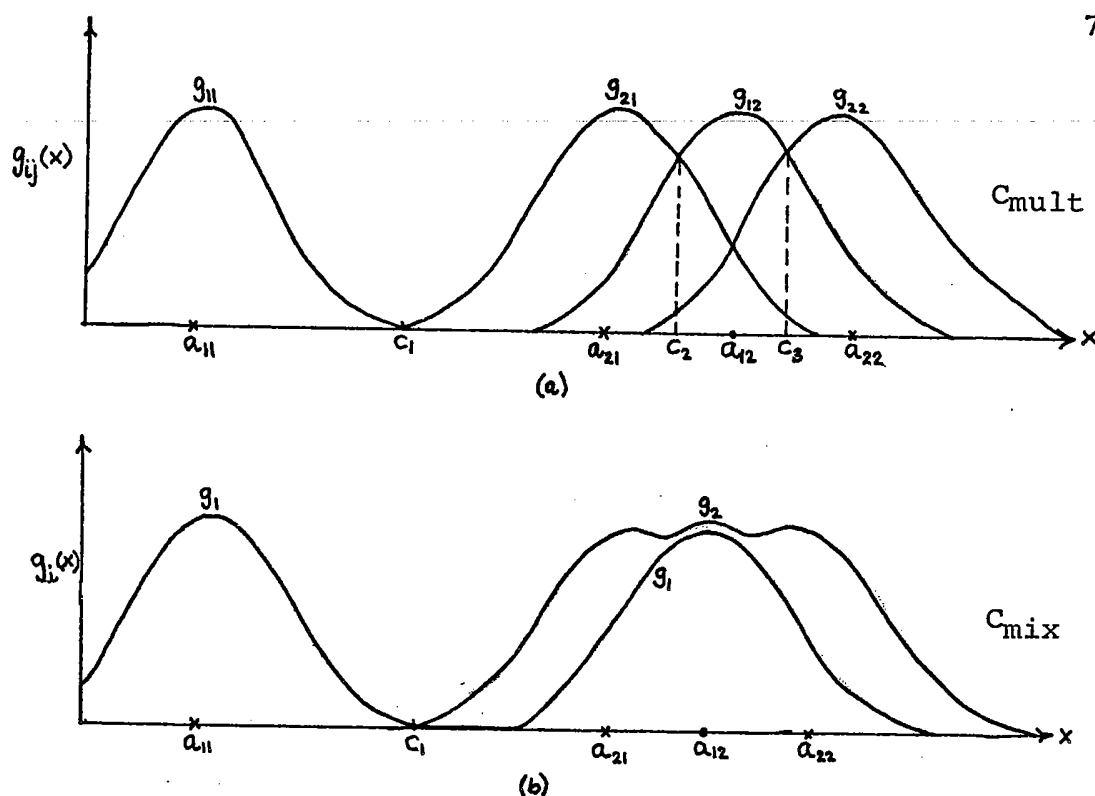


Fig. 4-2. Example of non-identity of C_{mult} and C_{mix}

The classifiers C_{mix} and C_{mult} provide a relationship between the Sebestyen and nearest-neighbor classifiers. If $p(\underline{x}/\omega_{ij})$ is multivariate Gaussian with mean \underline{a}_{ij} , and scalar covariance matrix $\sigma^2 I$, then C_{mix} is the Sebestyen classifier while C_{mult} is the nearest-neighbor classifier.

If $p(\underline{x}/\omega_{ij}) \sim n(\underline{a}_{ij}, \sigma^2 I)$, and the a priori probabilities $P(\omega_i)$, $P(\omega_{ij})$ are the same for every subclass, then

$$(4-1) \quad g_{ij}(\underline{x}) = \exp \left[-\|\underline{x} - \underline{a}_{ij}\|^2 / 2\sigma^2 \right].$$

The resulting decision surface of C_{mix} between classes ω_1 and ω_2 is

$$(4-2) \sum_{j=1}^S \left\{ \exp \left[-\frac{\|\underline{x} - \underline{a}_{1j}\|^2}{2\sigma^2} \right] - \exp \left[-\frac{\|\underline{x} - \underline{a}_{2j}\|^2}{2\sigma^2} \right] \right\} = 0 .$$

The decision surface defined by Eq. (4-2) is a continuous decision surface, which is also the decision surface of the Sebestyen classifier with $Q_j = (1/2\sigma^2)$, and $M_j = 1$ (see Section III.C).

The piecewise linear decision surface of C_{mult} is obtained from Eq. (4-1). It can be shown to have the form

$$(4-3) \quad \min_j \left\{ \|\underline{x} - \underline{a}_{1j}\|^2 \right\} - \min_j \left\{ \|\underline{x} - \underline{a}_{2j}\|^2 \right\} = 0 .$$

Eq. (4-3) defines the decision surface of the nearest-neighbor classifier based on the Euclidean distance measure. That is, if $p(\underline{x}/\omega_{ij}) \sim n(\underline{a}_{ij}, \sigma^2 I)$ and all a priori probabilities are equal, then C_{mix} and C_{mult} define the Sebestyen and nearest-neighbor classifiers respectively.

Having obtained a relationship between the Sebestyen and nearest-neighbor classifiers, we now proceed to make some observations concerning the decision surfaces of the two classifiers. The two decision surfaces are defined by Eqs. (4-2) and (4-3) in multi-dimensional feature-space. Since any point \underline{x} on the nearest-neighbor decision surface has to satisfy Eq. (4-3) it follows that \underline{x} is completely determined by a pair of reference patterns $(\underline{x}_1, \underline{x}_2)$ such

that $\underline{x}_1 \in \{a_{1j}\}$, $\underline{x}_2 \in \{a_{2m}\}$, and the following inequalities are satisfied

$$\begin{aligned} \|\underline{x} - \underline{x}_1\|^2 &\leq \|\underline{x} - a_{1j}\|^2 \\ \|\underline{x} - \underline{x}_2\|^2 &\leq \|\underline{x} - a_{2m}\|^2, \end{aligned}$$

for $j, m=1, \dots, s$. In other words, every point on the nearest-neighbor decision surface is completely determined by a pair of non-similarly classified reference patterns that are nearest to the point under consideration.

The decision boundaries of the Sebestyen and nearest-neighbor classifiers are illustrated in Fig. 4-3a for a particular case with $d=2$ and $s=2$. The piecewise linear decision boundary of the nearest-neighbor classifier is $c_1c_2c_3c_4$, where segment c_1c_2 is determined by the pair (a_{11}, a_{21}) , c_2c_3 is determined by the pair (a_{12}, a_{22}) , and c_3c_4 is determined by the pair (a_{12}, a_{22}) . The pair (a_{11}, a_{22}) does not contribute towards the determination of the nearest-neighbor decision boundary since any point \underline{x} in feature-space is equidistant from a_{11} and a_{22} is nearer to either a_{12} or a_{21} .

The decision boundary of the Sebestyen classifier in Fig. 4-3a is the broken line. It is a smoothed version of the piecewise linear decision boundary of the nearest-neighbor classifier. In the vicinity of point c_2 the Sebestyen decision boundary moves away from the nearest-

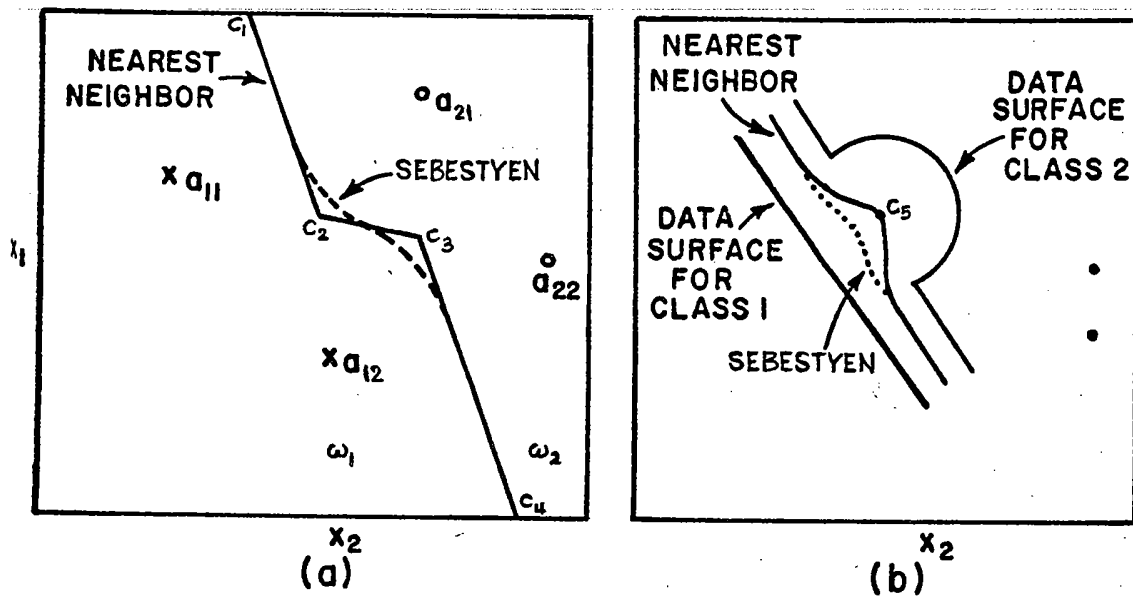


Fig. 4-3. Illustration of Sebestyen and nearest-neighbor decision surfaces.

nearest-neighbor decision boundary and towards the side of class ω_2 . This is explained by the fact that point c_2 is equidistant from points a_{11} , a_{21} , and a_{12} , while a_{22} is further away from c_2 than any of the other three points.

In the limit when s tends to infinity, and the data surfaces formed by ideal reference patterns are continuous, the decision surface of the nearest-neighbor classifier is also continuous. However, the resulting decision surface may not be identical to that of the Sebestyen classifier. This is illustrated in Fig. 4-3b where the ideal reference

pattern is considered to be a function of a single parameter, i.e., the ideal reference patterns lie on a curve in multivariate feature-space. In order to illustrate the fact that the Sebestyen and nearest-neighbor decision surfaces need not be identical we have chosen a special form of the curves corresponding to data surfaces. The data surface of class ω_2 is a straight line interrupted by a hump in the form of a semicircle (see Fig. 4-3b). The semicircle is centered at point c_5 . The data surface of class ω_1 is a straight line as shown. The continuous decision boundary of the nearest-neighbor classifier passes through point c_5 as depicted by the solid line. Since the summation corresponding to ω_1 exceeds the summation corresponding to class ω_2 for $\underline{x}=c_5$ in Eq. (4-2), point c_5 is not on the Sebestyen decision surface. Further, point c_5 is in that part of the feature-space associated with class ω_2 by the Sebestyen classifier. It follows that in the vicinity of c_5 , the Sebestyen decision surface is closer to the data surface of class ω_1 than the nearest-neighbor decision surface.

IV.C. Conditions for Similarity of the Sebestyen and Nearest-Neighbor Classifiers

When the location of reference patterns satisfy certain conditions in feature space, the Sebestyen classifier and the nearest-neighbor classifier based on the Euclidean distance measure have similar performance characteristics. We introduce here two notions of classifier similarity.

Identity: Two classifiers are identical for a specified reference set of patterns if their decision surfaces are identical.

ϵ -equivalence: Two classifiers are ϵ -equivalent if for specified class-conditional probability distributions, their average misclassification probabilities differ by no more than a specified positive ϵ .

Classifiers that are identical for certain reference pattern sets need not be identical for other reference pattern sets. In this section we will obtain conditions on reference pattern location for the identity of the Sebestyen and nearest-neighbor classifier.

IV.C.1 Conditions for the Identity of the Classifiers

Consider the case where each class is characterized by a single reference pattern. Let $\underline{a}_1 \in \omega_1$, and $\underline{b}_1 \in \omega_2$ be the

two reference patterns in d -space. The decision surfaces of the Sebestyen classifier and the nearest-neighbor classifier can be obtained from Eqs. (4-1) and (4-2), and shown to be identical to the hyperplane defined by

$$(\underline{a}_1 - \underline{b}_1)^t \underline{x} - \frac{1}{2}(\|\underline{a}_1\|^2 - \|\underline{b}_1\|^2) = 0.$$

That is, for $s=1$, the two classifiers are identical, irrespective of the location of the reference patterns.

In the case of $s=2$, consider the univariate problem with $a_1, a_2 \in \omega_1$, and $b_1, b_2 \in \omega_2$. The four points can occur in one of three possible distinct configurations on the real line. These are, writing sequentially in increasing order, $a_1 a_2 b_1 b_2$; $a_1 b_1 b_2 a_2$; and $a_1 b_1 a_2 b_2$. The other permutations are identical to one of the three except for the naming of points or classes. The configurations $a_1 b_1 b_2 a_1$ and $a_1 b_1 a_2 b_2$ imply that the data surfaces, which are the lines $a_1 a_2$ and $b_1 b_2$, intersect. This means that discrimination is impossible in region $b_1 b_2$ of the third configuration. We thus restrict our attention to the first configuration (see Fig. 4-4).

Irrespective of the distance of a_1 to a_2 and of b_2 to b_1 , the nearest-neighbor decision boundary is at $d_1 = (a_2 + b_1)/2$. That is, it classifies x as ω_1 if $x < d_1$, and

as ω_2 if $x > d_1$. However, the Sebestyen classifier takes into account

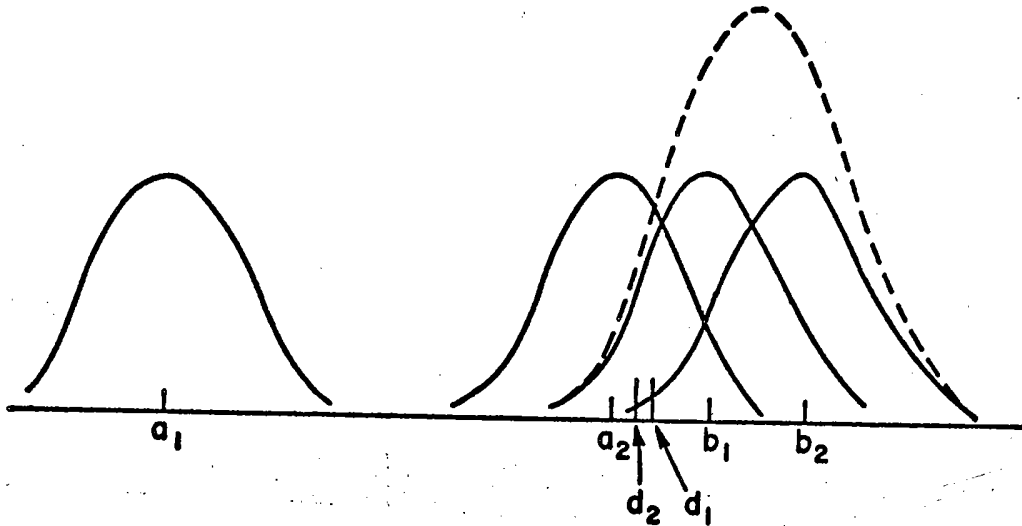


Fig. 4-4. Reference pattern configuration for obtaining conditions on classifier identity.

the weights provided by a_1 and b_2 . The decision boundary d_2 of the Sebestyen classifier is defined by the equation

$$(4-4) \quad \sum_{j=1}^2 \exp \left[-(d_2 - a_j)^2 / 2\sigma^2 \right] = \sum_{j=1}^2 \exp \left[-(d_2 - b_j)^2 / 2\sigma^2 \right].$$

Lemma 4.1: If $a_1, a_2 \in \omega_1$, $b_1, b_2 \in \omega_2$ such that $a_1 < a_2 < b_1 < b_2$, then the Sebestyen and nearest-neighbor classifiers are identical if and only if $(a_2 - a_1) = (b_2 - b_1)$.

Proof: Let d_1 and d_2 be the decision boundaries of the nearest-neighbor and Sebestyen classifiers. Let the classifiers be identical. That is $d_2 = d_1 = (a_2 + b_1)/2$.

Therefore,

$$\begin{aligned} & \exp\left[-((a_2 + b_1)/2 - a_1)^2 / 2\sigma^2\right] + \exp\left[-((b_1/2) - (a_2/2))^2 / 2\sigma^2\right] \\ &= \exp\left[-((a_2/2) - (b_1/2))^2 / 2\sigma^2\right] + \exp\left[-((a_2 + b_1)/2 - b_2)^2 / 2\sigma^2\right]. \end{aligned}$$

That is

$$\left[(a_2 + b_1)/2 - a_1\right]^2 = \left[(a_2 + b_1)/2 - b_2\right]^2.$$

It follows that either $a_1 = b_2$ or $(a_2 - a_1) = (b_2 - b_1)$. Since we have assumed $a_1 \neq b_2$, we have the necessary condition.

To prove sufficiency, let $(a_2 - a_1) = (b_2 - b_1)$. Thus the left-hand side of Eq. (4-4) is

$$\exp\left[-(d_2 + b_2 - b_1 - a_2)^2 / 2\sigma^2\right] + \exp\left[-(d_2 - a_2)^2 / 2\sigma^2\right]$$

or

$$\exp\left[-(d_2 - a_2)^2 / 2\sigma^2\right] \cdot \left\{ 1 + \exp\left[-(b_2 - b_1)^2 - 2(d_2 - a_2)(b_2 - b_1) / 2\sigma^2\right] \right\}.$$

Similarly, the right hand side of the equation can be written as

$$\exp\left[-(d_2 - b_1)^2 / 2\sigma^2\right] \cdot \left\{ 1 + \exp\left[-(a_2 - b_1)^2 - 2(d_2 - b_1)(a_1 - a_2) / 2\sigma^2\right] \right\}.$$

That is,

$$(d_2 - a_2)^2 = (d_2 - b_1)^2 \text{ or } d_2 = (a_2 + b_1)/2 = d_1. \quad \text{Q.E.D.}$$

A corollary of Lemma 4.1 provides the necessary and

sufficient condition for the identity of the Sebestyen and nearest-neighbor classifiers for the case where $s=2$, and $d=1$. The condition is that for each a_i there exists a point b_i such that a_i and b_i are at the same distance from the nearest-neighbor decision boundary d_1 . For $s>2$, the above condition on the a_i and b_i is sufficient for the identity of the two classifiers. However, it is not necessary, since the a_i 's and b_i 's can be situated on either side of d_1 such that the sum of their contributions toward the Sebestyen classifier decision boundary cancel at d_1 . The sufficient condition for the identity of the Sebestyen and nearest-neighbor classifiers can be extended to the general multivariate case as follows. For any point \underline{x} on the nearest-neighbor decision surface, and any positive scalar t , there exist an equal number of reference patterns of ω_1 and ω_2 that are at distance t from \underline{x} . This is stated in Theorem 4.1, the proof of which is given in Appendix B.

Theorem 4.1: Let $\{a_{1j}\} \in \omega_1$ and $\{a_{2j}\} \in \omega_2$ be ideal reference patterns. The nearest-neighbor and Sebestyen classifiers are identical if for every $\underline{\lambda}$ on the nearest-neighbor decision surface, and any a_{1j} distance t_j from $\underline{\lambda}$, there exists a unique a_{2j} distant t_j from $\underline{\lambda}$.

As a consequence of Theorem 4.1, the nearest-neighbor classifier will be identical to the Sebestyen classifier if the reference patterns are symmetric about the nearest-neighbor decision surface. Fig. 4-5a illustrates a particular configuration of points in 2-space for which the Sebestyen and nearest-neighbor decision surfaces are identical. Fig. 4-5b illustrates a case where the two classifiers are identical when the sufficient condition of Theorem 4.1 is not satisfied.

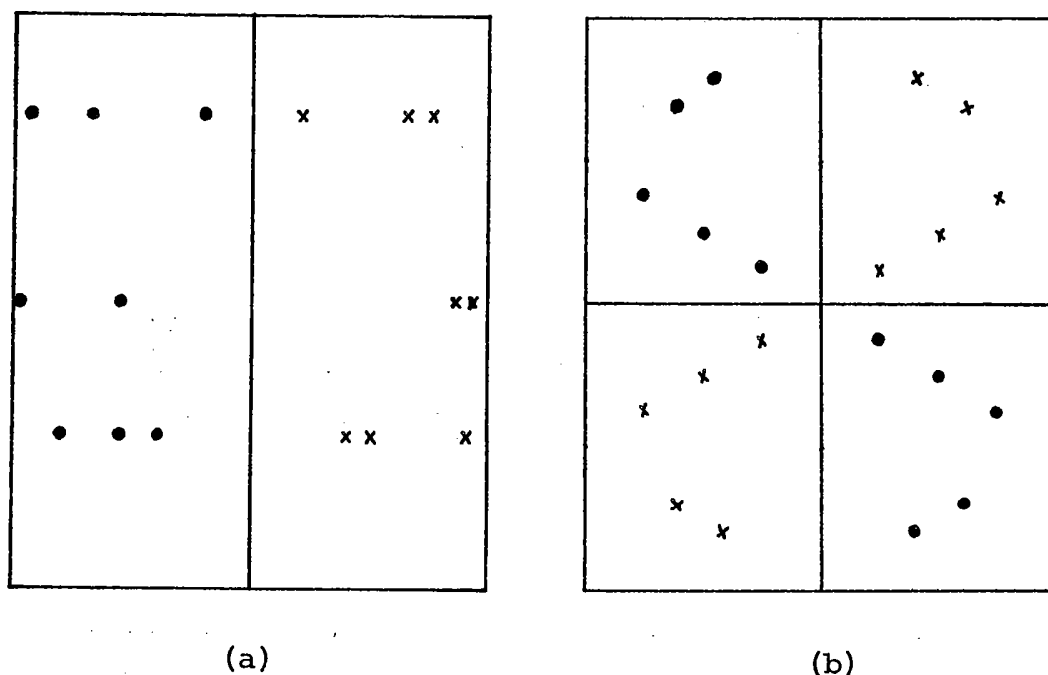


Fig. 4-5. Examples of Identity of the Sebestyen and Nearest-Neighbor Classifiers when $d=2$.

In the case where the nearest-neighbor classifier partitions the feature-space into two contiguous regions, necessary and sufficient conditions for classifier identity are provided by Theorem 4.2.

Theorem 4.2: If the nearest-neighbor classifier based reference points $\{\underline{a}_{1j}\} \in \omega_1$ and $\{\underline{a}_{2j}\} \in \omega_2$ is such that its decision surface, \mathcal{D} , partitions the feature-space into two contiguous regions, then the Sebestyen decision surface is identical to \mathcal{D} if and only if:

- (1) \mathcal{D} consists of a single hyperplane, and
- (2) for any reference point $\underline{\lambda}$,

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} \sum_{\underline{z} \in S_2(\underline{\lambda})} \left[\exp(-\|\underline{y} - \underline{z}\|^2 / 2\sigma^2) - \exp(-\|\underline{z} - \underline{z}\|^2 / 2\sigma^2) \right] = 0 ,$$

where $L(\underline{\lambda})$ = normal from $\underline{\lambda}$ to hyperplane \mathcal{D}

$$S_j(\underline{\lambda}) = \{ \underline{x} | \underline{x} \in \omega_j, \underline{x} \text{ lies on } L(\underline{\lambda}) \}$$

\underline{z} = point of intersection of $L(\underline{\lambda})$ and \mathcal{D} .

Proof: The first part of the proof is to show that the two conditions are necessary. We will first show that the nearest-neighbor decision surface consists of a single hyperplane. Since the two regions of feature-space are

contiguous it follows that \mathcal{D} has to be a connected surface. This implies that \mathcal{D} consists of intersecting subsets of hyperplanes. Such a decision surface cannot be represented by a continuous, differentiable function unless \mathcal{D} consists of a single hyperplane. Since the Sebestyen decision surface is represented by such a function it follows that \mathcal{D} has to be a single hyperplane.

We need to show that when the Sebestyen decision surface consists of a single hyperplane, the second condition of the theorem follows. When all reference points are located on a single line perpendicular to \mathcal{D} , the proof is trivial. Consider the case where all reference points lie on two lines that are perpendicular to \mathcal{D} . Let $\underline{\lambda}_1$ and $\underline{\lambda}_2$ be two reference points located on each of the two lines. Let $\underline{\xi}_1$ and $\underline{\xi}_2$ be the corresponding points of intersection of the two lines with \mathcal{D} .

Abbreviating $\exp(-\|\alpha - \beta\|^2/2\sigma^2)$ as $f(\alpha, \beta)$, we can write two equations corresponding to each of $\underline{\xi}_1$ and $\underline{\xi}_2$ as

$$\sum_{i=1}^2 \sum_{\underline{y} \in S_1(\underline{\lambda}_i)} \sum_{\underline{z} \in S_2(\underline{\lambda}_i)} \left[f(\underline{\xi}_j, \underline{y}) - f(\underline{\xi}_j, \underline{z}) \right] = 0$$

where $j=1,2$. For any $\underline{x} \in S_k(\underline{\lambda}_i)$, we can write

$$\|\underline{x} - \underline{\xi}_j\|^2 = \|\underline{\xi}_1 - \underline{\xi}_2\|^2 + \|\underline{x} - \underline{\xi}_i\|^2$$

for $i, j = 1, 2$, $i \neq j$, and $k=1, 2$. Thus the two equations

corresponding to ξ_1 and ξ_2 are of the form

$$A_1 + c_{12}A_2 = 0$$

$$c_{12}A_1 + A_2 = 0 ,$$

where

$$A_j = \sum_{y \in S_1(\lambda_j)} \sum_{z \in S_2(\lambda_j)} [f(\xi_j, y) - f(\xi_j, z)] \quad j=1,2 ,$$

and $c_{12} = \exp(-\|\xi_1 - \xi_2\|^2 / 2\sigma^2)$. The value of c_{12} is dependent on the distance between the two lines, which in turn is an arbitrary value. It follows that the two equations are simultaneously satisfied if and only if $A_1 = A_2 = 0$.

In general all reference points will be located on m lines perpendicular to \mathcal{D} , for some finite integer m . In such a case we will have a set of m equations of the form

$$A_j + \sum_{\substack{k=1 \\ k \neq j}}^m c_{jk} A_k = 0 , \text{ for } j=1, \dots, m.$$

For arbitrary values of c_{jk} , these equations are simultaneously satisfied if and only if $A_j = 0$ for all j . This completes the proof that conditions (1) and (2) are necessary.

We need to show that when the nearest-neighbor decision surface is a single hyperplane \mathcal{D} and condition (2)

corresponding to \underline{x}_1 and \underline{x}_2 are of the form

$$A_1 + c_{12}A_2 = 0$$

$$c_{12}A_1 + A_2 = 0 ,$$

where

$$A_j = \sum_{\underline{y} \in S_1(\underline{x}_j)} \sum_{\underline{z} \in S_2(\underline{x}_j)} [f(\underline{x}_j, \underline{y}) - f(\underline{x}_j, \underline{z})] \quad j=1,2 ,$$

and $c_{12} = \exp(-\|\underline{x}_1 - \underline{x}_2\|^2 / 2\sigma^2)$. The value of c_{12} is dependent on the distance between the two lines, which in turn is an arbitrary value. It follows that the two equations are simultaneously satisfied if and only if $A_1 = A_2 = 0$.

In general all reference points will be located on m lines perpendicular to \mathcal{D} , for some finite integer m . In such a case we will have a set of m equations of the form

$$A_j + \sum_{\substack{k=1 \\ k \neq j}}^m c_{jk} \cdot A_k = 0 , \text{ for } j=1, \dots, m.$$

For arbitrary values of c_{jk} , these equations are simultaneously satisfied if and only if $A_j = 0$ for all j . This completes the proof that conditions (1) and (2) are necessary.

We need to show that when the nearest-neighbor decision surface is a single hyperplane \mathcal{D} and condition (2)

of the theorem is satisfied, the Sebestyen decision surface is identical to \mathcal{D} . This will be shown by first considering an arbitrary point $\underline{\delta}_1$ on \mathcal{D} and showing that it has to be a point on the Sebestyen decision surface. We will then consider a point $\underline{\delta}_2$ that does not lie on \mathcal{D} and show that the point cannot lie on the Sebestyen decision surface.

Let $\underline{\delta}_1$ be an arbitrary point on \mathcal{D} . According to condition (2), for any reference point $\underline{\lambda}$, we can write

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} \sum_{\underline{z} \in S_2(\underline{\lambda})} [f(\underline{y}, \underline{\xi}) - f(\underline{z}, \underline{\xi})] = 0$$

where $\underline{\xi}$ is defined by $\underline{\lambda}$ and \mathcal{D} . Since for $\underline{x} \in S_j(\underline{\lambda})$, $j=1,2$, we can write $\|\underline{x} - \underline{\delta}_1\|^2 = \|\underline{x} - \underline{\xi}\|^2 + \|\underline{\xi} - \underline{\delta}_1\|^2$, it follows that for any reference point $\underline{\lambda}$,

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} \sum_{\underline{z} \in S_2(\underline{\lambda})} [f(\underline{y}, \underline{\delta}_1) - f(\underline{z}, \underline{\delta}_1)] = 0.$$

By considering equations corresponding to every reference point, we can write

$$\sum_{\underline{y} \in \omega_1} \sum_{\underline{z} \in \omega_2} [f(\underline{y}, \underline{\delta}_1) - f(\underline{z}, \underline{\delta}_1)] = 0.$$

It follows that $\underline{\delta}_1$ is a point on the Sebestyen decision surface.

Let $\underline{\delta}_2$ be an arbitrary point on the ω_1 side of hyperplane \mathcal{D} . Let $\underline{\lambda}$ be an arbitrary reference point such that $\underline{\lambda} \neq \underline{\delta}_2$. Let \underline{y} and \underline{z} be the points of intersection \mathcal{D} with $L(\underline{\delta}_2)$ and $L(\underline{\lambda})$ respectively. For any $\underline{y} \in S_1(\underline{\lambda})$, and $\underline{z} \in S_2(\underline{\lambda})$, we can write

$$\begin{aligned}\|\underline{y} - \underline{\delta}_2\|^2 &\leq \|\underline{y} - \underline{r}\|^2 + \|\underline{r} - \underline{\delta}_2\|^2 \leq \|\underline{y} - \underline{x}\|^2 + \|\underline{x} - \underline{r}\|^2 + \|\underline{r} - \underline{\delta}_2\|^2 \\ \|\underline{z} - \underline{\delta}_2\|^2 &\geq \|\underline{z} - \underline{r}\|^2 + \|\underline{r} - \underline{\delta}_2\|^2 \geq \|\underline{z} - \underline{x}\|^2 + \|\underline{x} - \underline{r}\|^2 + \|\underline{r} - \underline{\delta}_2\|^2.\end{aligned}$$

It follows that

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} f(\underline{y}, \underline{\delta}_2) \geq K \cdot \sum_{\underline{y} \in S_1(\underline{\lambda})} f(\underline{y}, \underline{x})$$

and,

$$\sum_{\underline{z} \in S_2(\underline{\lambda})} f(\underline{z}, \underline{\delta}_2) \leq K \cdot \sum_{\underline{z} \in S_2(\underline{\lambda})} f(\underline{z}, \underline{x}),$$

where $K = f(\underline{x}, \underline{r}) + f(\underline{r}, \underline{\delta}_2)$. According to the second condition of the theorem

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} f(\underline{y}, \underline{x}) = \sum_{\underline{z} \in S_2(\underline{\lambda})} f(\underline{z}, \underline{x}).$$

Thus we have the inequality

$$\sum_{\underline{y} \in S_1(\underline{\lambda})} f(\underline{y}, \underline{\delta}_2) \geq \sum_{\underline{z} \in S_2(\underline{\lambda})} f(\underline{z}, \underline{\delta}_2).$$

Since the last inequality is true for every reference point $\underline{\lambda}$, it follows that the Sebestyen classifier considers $\underline{\delta}_2$ to be a point from ω_1 . Similarly we can show that if $\underline{\delta}_2$ were a point on the ω_2 side of \mathcal{D} , the Sebestyen classifier would classify it as a point from ω_2 . We have thus shown that condition (1) and (2) are sufficient for classifier identity.

Q.E.D.

IV.C.2 Algorithm for Determining Classifier Identity

Whenever the Sebestyen and nearest-neighbor classifiers are identical, their decision surface consists of a set of parallel or orthogonal hyperplanes. The following algorithm, called Algorithm C, determines identity when the decision surface consists of a single hyperplane. With some simple modifications, the algorithm can be extended to the situation where decision surfaces can consist of more than one hyperplane. The only inputs to Algorithm C are the two sets of reference points $\{a_{1i}\}$ and $\{a_{2j}\}$, which represent classes ω_1 and ω_2 respectively.

The algorithm begins by determining the nearest pair of dissimilarly classified reference points. A hyperplane determined by this pair is used as a candidate for the overall decision surface. Subsequent tests for classifier identity are based on conditions provided by Theorems 4.1 and 4.2.

Algorithm C: Classifier Identity

1. Determine the nearest pair of points a_{1m} and a_{2n} .
2. Determine hyperplane H , the perpendicular bisector of line segment joining a_{1m} and a_{2n} .

If H does not separate classes ω_1 and ω_2 , identity

cannot be determined. However, the algorithm can be extended to resolve this case.

3. Let M be a binary relation on $U = \{\underline{a}_{1i}\} \cup \{\underline{a}_{2j}\}$,

defining matched pairs of points, as follows:

$$M = \left[(\underline{a}_{1i}, \underline{a}_{2j}) \mid \begin{array}{l} H \text{ is the perpendicular bisector of a} \\ \text{line segment joining } \underline{a}_{1i} \text{ and } \underline{a}_{2j} \end{array} \right]$$

Let (S_1, S_2) be a partition of U as follows:

$$S_1 = [\underline{x}, \underline{y} \mid \underline{x}, \underline{y} \in U, \text{ and } \underline{xMy}]$$

$$S_2 = U - S_1 .$$

Determine S_1 and S_2 .

4. If $S_2 = \emptyset$, the two classifiers are identical as a consequence of Theorem 4.1.

5. Determine a hyperplane H_1 , defined as follows:

H_1 is orthogonal to H , passes through a pair of points $(\underline{x}, \underline{y})$ such that \underline{xMy} , and all points in U lie in the same closed half space of H_1 .

6. Initialize for loop:

$$\text{let } k=1, \text{ and } C_1 = C'_1 = [\underline{x} \mid \underline{x} \in U, \text{ and } \underline{x} \text{ lies on } H_1] .$$

7. Let $\underline{\lambda} \in C_k$. Determine the set $S_{\underline{\lambda}}$, defined as follows:

$$S_{\underline{\lambda}} = [\underline{x} \mid \underline{x} \in C_k, \text{ and } \underline{x} \text{ lies on } L(\underline{\lambda})] , \text{ where } L(\underline{\lambda})$$

is defined by $\underline{\lambda}$ and hyperplane H .

8. If the following exponential constraint is not

satisfied the classifiers are not identical due to
to Theorem 4.2:

$$\sum_{\underline{y} \in \omega_1 \cap S_\lambda} \sum_{\underline{z} \in \omega_2 \cap S_\lambda} \exp\left[-\|\underline{x} - \underline{y}\|^2 / 2\sigma^2\right] - \exp\left[-\|\underline{x} - \underline{z}\|^2 / 2\sigma^2\right] = 0 .$$

9. Delete S_λ from C_k . If $C_k \neq \emptyset$, go to 7.
10. Increment k by one. Determine set C_k , defined as

$$C_k = \left[\underline{x} \mid \begin{array}{l} \underline{x} \in \{U - C'_{k-1} - \dots - C'_1\} , \text{ and} \\ d(\underline{x}, H_{k-1}) = \min_{\underline{x} \in U} d(\underline{x}, H_{k-1}) \end{array} \right] .$$

Let $C'_k = C_k$.

11. If $C_k = \emptyset$, go to 13 (at this point the Sebestyen decision surface is the hyperplane H).
12. Determine hyperplane H_k defined as follows:

H_k passes through all points in C_k ,

and is parallel to H_1, \dots, H_{k-1} .

Go to 7.

13. Determine whether the nearest-neighbor decision surface is the hyperplane H .

If so, classifiers are identical; otherwise, they are not identical.

With the exception of the test made in Step 13, all steps in Algorithm C can be implemented in a straight-forward manner. The test in Step 13 is a special case of the general problem of

determining the nearest-neighbor decision surface. This is a complex problem which has been discussed recently in the literature. (see Dasarathy, 1975)

In order to obtain an algorithm for cases where the decision surface is not a single hyperplane, we proceed as follows. Let H partition the feature-space into regions h_1 and h_2 . Let there be points of more than one class in region h_1 . Determine the hyperplane K that separates the closest pair of dissimilarly classified points in h_1 . If K is neither orthogonal nor parallel to H , classifiers are not identical. If either condition is true, let (k_1, k_2) be the feature-space partition due to K . Determine if the partition $(h_1 \cap k_1, h_1 \cap k_2, h_2 \cap k_1, h_2 \cap k_2)$ is such that each block of the partition contains no pair of points of dissimilar classification. If there exists such a pair, refine the partition using a hyperplane L that is either parallel or orthogonal to H and K . Continue the process of partition refinement until all points of a given block are similarly classified.

Let H, K, L, \dots, Z be n candidate hyperplanes obtained as a result of the above procedure. We now apply Algorithm C beginning at Step 3, to each of the n hyperplanes. In each application, the candidate hyperplane is considered in place of H . The classifiers are identical if and only if all n hyperplanes satisfy the test for identity.

IV.D. Robustness of the Sebestyen and Nearest-Neighbor Classifiers

The decision surfaces of the Sebestyen and nearest-neighbor classifiers are dependent upon the values of stored parameters. The only stored parameters of the nearest-neighbor classifier are ideal reference patterns, whereas the parameters of the Sebestyen classifier also include the noise standard deviation, σ . The average misclassification probabilities of the two classifiers can be expected to be lowest when stored parameters have the same values as the true parameters of class-conditional distributions they represent. In this section the dependence of the Sebestyen decision boundary on the stored value of σ is examined. An analysis is made in the univariate case, and from the point of view of inexact knowledge of the true value of σ .

It was observed in the previous section that when reference patterns satisfy certain location constraints the Sebestyen and nearest-neighbor decision surfaces are identical. These constraints, which specify reference pattern location, do not include the value of σ . That is, the performance of the Sebestyen classifier is independent of the value of σ when such conditions on reference pattern location are satisfied. We will first consider a reference

pattern configuration that leads to different decision boundaries of the two classifiers. Some guidelines for the assumption of σ will be obtained. This is followed by a consideration of the general univariate problem. An expression for the rate of variation of the Sebestyen decision boundary with the assumed value of σ will be obtained. Based on this expression it can be concluded that an underestimation of σ is preferable to overestimation when σ values lie in the range of values of interest in the radar aircraft identification problem.

Consider the linearly separable, univariate problem with two reference patterns in each class. Let $a_1, a_2 \in \omega_1$, and let $b_1, b_2 \in \omega_2$. Further let $a_1 = a_2$, and $b_2 \gg b_1 > a_1$. This reflects a situation where the reference patterns of ω_1 are clustered together, as compared to the reference patterns of ω_2 (Fig. 4-6). If σ is the assumed value of noise standard deviation, then the decision boundary of the Sebestyen classifier, η , is defined by

$$2 \cdot \exp\left[-(\eta - a_1)^2 / 2\sigma^2\right] = \exp\left[-(\eta - b_1)^2 / 2\sigma^2\right].$$

Solving for η , we have the Sebestyen decision boundary as a function of σ as follows

$$(4-5) \quad \eta = \eta(\sigma) = (b_1 + a_1) / 2 + \sigma^2 \ln 2 / (b_1 - a_1).$$

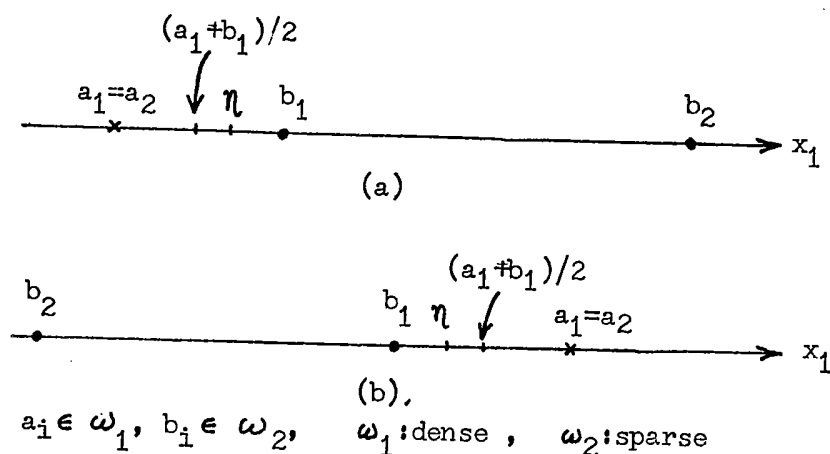


Fig. 4-6. Reference pattern configuration to illustrate dependence of the Sebestyen decision boundary on the assumed value of σ .

That is, the decision boundary is a quadratic in σ , and monotone increasing with σ ; the latter is observed by the fact that the first derivative of η with respect to σ is positive. Since the decision boundary lies in the interval (a_1, b_1) it follows that as the classifier assumption of σ increases, the decision boundary moves towards ω_2 .

Consider now the problem situation where the true value of noise standard deviation, σ_t , is not known exactly. If σ_t is known to be uniformly distributed in the interval (σ_1, σ_2) then the following questions arise: What should be the assumed value of σ , and which of the assumptions

$\sigma = \sigma_1$ or $\sigma = \sigma_2$ is preferable? The answers to these questions are provided by Theorem 4.3, the proof of which is based on the following lemma.

Lemma 4.3: If σ is uniformly distributed in the interval (a, b) , $f(\sigma)$ is a quadratic in σ , and $E[f(\sigma)]$ is the expected value of $f(\sigma)$, then

$$E[f(\sigma)] - f(a) \leq f(b) - E[f(\sigma)] .$$

Proof: Let $f(\sigma) = a_0\sigma^2 + a_1\sigma + a_2$. Therefore,

$$\begin{aligned} E[f(\sigma)] &= \frac{1}{b-a} \int_a^b f(\sigma) \cdot d\sigma \\ &= \frac{1}{b-a} \left[\frac{a_0}{3} \cdot (b^3 - a^3) + \frac{a_1}{2} \cdot (b^2 - a^2) + a_2(b-a) \right] , \end{aligned}$$

and

$$f(a) = a_0a^2 + a_1a + a_2 .$$

$$\text{That is, } E[f(\sigma)] - f(a) = (b-a) \left[\frac{a_0}{3}(b+2a) + a_1/2 \right] .$$

Similarly we can show that

$$f(b) - E[f(\sigma)] = (b-a) \left[\frac{a_0}{3}(a+2b) + a_1/2 \right] .$$

Since $b > a$, the result follows.

Q.E.D.

The result can be shown to be true when $f(\sigma)$ is a higher order polynomial; a quadratic suffices for our purpose.

However, the result is not true for transcendental functions of the form e^σ . Theorem 4.3 assumes that the optimal

decision boundary is one that is equal to the expected value of $\eta(\sigma_t)$, where σ_t is the true noise standard deviation. The expected value is obtained with respect to the distribution of σ_t .

Theorem 4.3: Let $a_1, a_2 \in \omega_1$, $b_1, b_2 \in \omega_2$ be four equally likely subclass means such that $a_1 = a_2$, $b_2 > b_1 > a_1$. Let $n(\mu, \sigma_t^2)$ be the distribution of the noisy pattern when μ is the true subclass mean. If σ_t is uniformly distributed in the interval (σ_1, σ_2) , then an assumption of $\sigma_{\text{opt}} = [(\sigma_2^2 + \sigma_1 \cdot \sigma_2 + \sigma_1^2)/3]^{1/2}$ by the Sebestyen classifier results in the optimal decision boundary. An assumption of σ_1 leads to a smaller change from the optimal decision boundary than an assumption of σ_2 .

Proof: Since σ_t is uniform in the interval (σ_1, σ_2) , and the decision boundary is as defined by Eq. (4-7), it

follows that

$$\begin{aligned} E(\eta(\sigma_t)) &= \int_{\sigma_1}^{\sigma_2} [1/(\sigma_2 - \sigma_1)] \cdot \left\{ (b_1 + a_1)/2 + [\sigma^2 \ln 2 / (b_1 - a_1)] \right\} \cdot d\sigma \\ &= (b_1 + a_1)/2 + [\ln 2 \cdot (\sigma_2^3 - \sigma_1^3) / [3 \cdot (b_1 - a_1) \cdot (\sigma_2 - \sigma_1)]] \end{aligned}$$

Solving $E(\eta(\sigma_t)) - \eta(\sigma_{\text{opt}}) = 0$, we get the value of σ_{opt} .

The second half of the theorem follows from Lemma 4.3.

That is, since η is a quadratic in σ ,

$$E[\eta(\sigma_t)] - \eta(\sigma_1) \leq \eta(\sigma_2) - E[\eta(\sigma_t)] \quad \text{Q.E.D.}$$

If we assumed $a_1 > b_1 \gg b_2$, i.e., interchanged the relative positions of the dense and sparse clusters as shown in Fig. 4-6b, η would be monotone decreasing with σ . However, the decision boundary continues to move towards the sparse cluster with increasing values of σ ; σ_1 can be shown to be a better assumption than σ_2 by showing that

$$\eta(\sigma_1) - E(\eta(\sigma_t)) < E(\eta(\sigma_t)) - \eta(\sigma_2).$$

In a particular problem situation if one is interested in knowing if $\sigma_t + \delta$ is a better assumption than $\sigma_t - \delta$, when σ_t is the true noise standard deviation, then the functional forms of η and its derivative with respect to σ provide the answer. In order to show this, we need the following lemma.

Lemma 4.4: If $f(\sigma)$ and its derivative $f'(\sigma)$ are continuous and monotone increasing with σ in a closed interval I , then for any $\sigma_1, \sigma_1 \pm \delta \in I$,

$$f(\sigma_1) - f(\sigma_1 - \delta) \leq f(\sigma_1 + \delta) - f(\sigma_1).$$

Proof: From the mean value theorem for derivatives we can write

$$f(\sigma_1) - f(\sigma_1 - \delta) = f'(\xi_1) \cdot \delta \quad \text{where } \xi_1 \in (\sigma_1 - \delta, \sigma_1)$$

$$f(\sigma_1 + \delta) - f(\sigma_1) = f'(\xi_2) \cdot \delta \quad \text{where } \xi_2 \in (\sigma_1, \sigma_1 + \delta).$$

Since $\xi_2 \geq \xi_1$, the monotonicity of $f'(\sigma)$ implies that $f'(\xi_2) \geq f'(\xi_1)$, and hence the result. Q.E.D.

It can be seen from Eq. (4-5) that the first and second derivatives of η with respect to σ are both positive for all positive σ . It follows from Lemma 4.4 that

$$|\eta(\sigma_t - \delta) - \eta(\sigma_t)| \leq |\eta(\sigma_t + \delta) - \eta(\sigma_t)|,$$

or that $\sigma_t - \delta$ leads to a smaller change from the optimal decision boundary than $\sigma_t + \delta$. That is, we do better by underestimating the value of σ_t than by overestimating it. Again, this conclusion is independent of the relative class location.

In the more complex situation where each class is characterized by more than two reference patterns, η cannot be expressed explicitly as a function of σ . However, intuition into the behavior of η with σ can be obtained by examining the derivative of η with respect to σ . If $a_{1i} \in \omega_1$ and $a_{2i} \in \omega_2$, the decision boundary, η , of the Sebestyen classifier is defined implicitly by the equation

$$F(\eta, \sigma) = \sum_{i=1}^S \left[\exp(-A_{1i}^2 / 2\sigma^2) - \exp(-A_{2i}^2 / 2\sigma^2) \right] = 0$$

where $A_{ji} = (\eta - a_{ji})$. The derivative of η with respect to

σ is obtained from the identity

$$\frac{\partial F}{\partial \sigma} + \frac{\partial F}{\partial \eta} \frac{d(\eta)}{d\sigma} = 0 .$$

That is,
$$\frac{d(\eta)}{d\sigma} = \frac{-\partial F / \partial \sigma}{\partial F / \partial \eta} .$$

Now,

$$\begin{aligned} \frac{\partial F}{\partial \sigma} &= (2/\sigma^3) \sum_{i=1}^S A_{1i}^2 \exp(-A_{1i}^2/2\sigma^2) - A_{2i}^2 \exp(-A_{2i}^2/2\sigma^2) \\ \frac{\partial F}{\partial \eta} &= -(1/\sigma^2) \sum_{i=1}^S A_{1i} \exp(-A_{1i}^2/2\sigma^2) - A_{2i} \exp(-A_{2i}^2/2\sigma^2) . \end{aligned}$$

That is,

$$(4-6) \quad \frac{d(\eta)}{d\sigma} = (2/\sigma) \frac{\sum_{i=1}^S A_{1i}^2 \exp(-A_{1i}^2/2\sigma^2) - A_{2i}^2 \exp(-A_{2i}^2/2\sigma^2)}{\sum_{i=1}^S A_{1i} \exp(-A_{1i}^2/2\sigma^2) - A_{2i} \exp(-A_{2i}^2/2\sigma^2)} .$$

Eq. (4-6) is an explicit representation of the rate of change of the decision boundary of the Sebestyen classifier, with the value of noise standard deviation, in the general univariate case. Considering the situation where each class is represented by an arbitrary number of reference patterns, Eq. (4-6) can be used to show that if one of the classes is sparse with respect to the other, an underestimation of σ is preferable to an overestimation. In the following

discussion we will consider the linearly separable univariate problem and show that an underestimation of σ is preferable to an overestimation. It is assumed that one of the classes is sparse with respect to the other, and the true and estimated values of σ lie in a specified interval.

If ω_1 and ω_2 are linearly separable with $a_{1i} < a_{2i}$, i.e., $A_{1i} > A_{2i}$, $d(\eta)/d\sigma$ can be verified to be positive if $\sigma \in (\frac{1}{2}f_2, f_1)$ and negative if $\sigma \in (\frac{1}{2}f_1, f_2)$ where $f_1 = \max_i f_i(\eta)$, $f_2 = \min_i f_i(\eta)$, and $f_i(\eta) = A_{1i}^2 - A_{2i}^2 / \ln(A_{1i}^2 / A_{2i}^2)$. When $A_{1i}^2 = A_{2i}^2$, we can show that $f_i(\eta) = A_{1i}^2 = A_{2i}^2$ by the following limiting argument

$$\lim_{x \rightarrow a} \left\{ \frac{a-x}{\ln(a/x)} \right\} = \lim_{x \rightarrow a} \left\{ \frac{-1}{-1/x} \right\} = a.$$

It follows that if η is at the mid-point between a_{1i} and a_{2i} , then $f_i(\eta) = \delta^2/4$, where $\delta = |a_{1i} - a_{2i}|$. If the closest points of each class are separated by a distance λ , which is the inter-class separation, then over all possible naming of points, at least for one value of i , $f_i(\sigma)$ has the value $\lambda^2/4$.

In the particular case where the reference patterns of ω_1 are dense with respect to those of ω_2 , i.e., the average intra-class separation of ω_1 is small compared to the average intra-class separation of ω_2 , values of A_{1i} tend to be much larger than the values A_{2i} . For large values of A_{1i} as compared to A_{2i} , $f_i(\sigma)$ can be approximated by

$$g_i(\eta) = A_{1i}^2 / [\ln A_{1i}^2],$$

which follows from the fact that the relative error, defined by

$$\frac{f_i(\eta) - g_i(\eta)}{f(\eta)} = \frac{A_{1i}^2 \ln A_{2i}^2 - A_{2i}^2 \ln A_{1i}^2}{(A_{1i}^2 - A_{2i}^2) \ln A_{1i}^2}$$

is a decreasing function of A_{1i} for constant A_{2i} . The minimum value of $g_i(\eta)$ over all possible value of A_{1i} is e . Thus, if the points of ω_2 are all closely clustered together, and the points of ω_1 are sparse, $d(\eta)/d\sigma$ is positive in the range of $(0, \frac{1}{2}f_2) = [0, \frac{1}{2}\min(\lambda^2/4, e)]$. If the inter-class separation, λ , is greater than $2e^{\frac{1}{2}}$, then $\min(\lambda^2/4, e) = e$, and η is guaranteed to move towards the sparse cluster as long as $\sigma^2 < e^2/4 = 1.82$. In the case where the inter-class separation is significantly less than

$2e^{\frac{1}{2}}$, a transformed feature space can be found so that $\min(\lambda^2/4, e) = e$.

The second derivative $d^2(\eta)/d\sigma^2$ can be obtained by direct differentiation of the above expression for $d(\eta)/d\sigma$. The sign of $d^2(\eta)/d\sigma^2$ is determined by the sign of the expression

$$\sum_{i=1}^S A'_{1i} \exp[-A_{1i}^2/2\sigma^2] + A'_{2i} \exp[-A_{2i}^2/2\sigma^2] + A' \exp[-(A_{1i}^2 + A_{2i}^2)/2\sigma^2]$$

where $A'_{ji} = A_{ji}^5 (4/\sigma^2 - 2/\sigma^3) + 2A_{ji}^3$

$$A' = 2 \cdot \left[(2/\sigma^3) (A_{1i}^4 A_{2i} + A_{2i}^4 A_{1i}) - (2/\sigma^2) (A_{1i}^3 A_{2i}^2 + A_{2i}^3 A_{1i}^2) + (A_{1i}^2 A_{2i}^2 + A_{2i}^2 A_{1i}^2) \right]$$

Since we have assumed $A_{1i} \gg A_{2i}$, the sign of each term of the summation is determined primarily by the term

$$A_{1i}^5 (4/\sigma^2 - 2/\sigma^3) + 2A_{1i}^3 \exp(-A_{1i}^2/2\sigma^2)$$

Now if $(4/\sigma^2 - 2/\sigma^3) > 0$, every one of the above terms will be positive. That is, if $\sigma > \frac{1}{2}$, or $\sigma^2 > \frac{1}{4}$, $d^2(\eta)/d\sigma^2$ is guaranteed to be positive.

We have thus been able to show that when one of the classes is densely clustered with respect to the other, and the classes are linearly separable, then the Sebestyen decision boundary and its derivative with respect to σ are both monotone increasing in a range of values that are of interest in

the radar aircraft identification problem. It follows from Lemma 4.5 that if σ_t is the true noise standard deviation, then $\sigma_t - \delta$ is a better assumption than $\sigma_t + \delta$ in the above range of values. This conclusion again pertains to the linearly separable univariate problem. Further, a small change in the decision boundary from its optimal value can be expected to lead to a correspondingly small change in the misclassification probability. Experimental results in the multivariate case with radar aircraft data indicate that an underestimation of σ leads to a smaller increase from the minimum achievable misclassification probability than an overestimation of σ (see Section V.D.1).

IV.E. Effect of Increasing Dimensionality on the Sebestyen and Nearest-Neighbor Classifiers

In this section the effect of increasing the number of features on the performance of the Sebestyen and nearest-neighbor classifiers is examined. When multiple orientation possibilities are considered and dimensionality is increased, certain counter-intuitive phenomena can occur, which will be considered here. We start by demonstrating that the average misclassification probabilities of each of the two classifiers has to decrease with dimensionality, as would be intuitively

expected.

If class ω_i is characterized by the set of reference patterns \underline{a}_{ij} , $j=1, \dots, s$ and noisy patterns corresponding to \underline{a}_{ij} are distributed according to $n(\underline{a}_{ij}, \sigma^2 \mathbf{I})$, then the average misclassification probability of the Sebestyen classifier can be approximated by the expression obtained in Section III.D., i.e.,

$$P_e = (1/2s^2) \sum_{j=1}^s \sum_{k=1}^s \exp \left[-\|\underline{a}_{1j} - \underline{a}_{2k}\|^2 / 8\sigma^2 \right] .$$

The Euclidean norm $\|\underline{a}_{1j} - \underline{a}_{2k}\|$ is a non-decreasing function of the dimensionality, d , of the feature-space for all values of the components of the reference patterns \underline{a}_{ij} . Thus the average misclassification probabilities of the Sebestyen and nearest-neighbor classifiers are non-increasing functions of d .

In spite of the fact that the average misclassification probability between classes decreases, the misclassification probability between classes at particular orientation possibilities can actually increase with the introduction of new features. This phenomenon can occur with either the nearest-neighbor or Sebestyen classifier, although the causes for the phenomenon need not be the same for the two classifiers. In

the rest of this section we obtain sufficient conditions for observing a direct relationship between misclassification probability and dimensionality. The conditions indicate the type of reference pattern configuration and the type of feature that leads to the phenomenon. The knowledge of such conditions is useful in the following problem situation: a classifier implemented on the basis of multiple orientation possibilities is to be utilized in a case where the orientations with non-zero probability form a proper subset of the orientations considered by the classifier. For example, a classifier that has been designed on the basis of all aircraft orientation possibilities may have to be utilized in a situation where the aircraft orientation is known approximately. Under this problem constraint we can ask the question whether the average misclassification probability for particular orientations, i.e., between subclasses corresponding to particular orientations, will decrease with the incorporation of a new feature into the classifier. The answer to this query depends upon the relative movement of the reference points with the change in dimensionality.

IV.E.1 Dimensionality and the Sebestyen Classifier

Consider the use of the Sebestyen classifier to distinguish between two classes that have s subclasses in each class. We will show that if the classifier utilizes d features, and a $(d+1)^{\text{th}}$ feature is added to the classifier, the average misclassification probability between the q^{th} subclass of each class will increase with the transition in feature-space from E^d to E^{d+1} when the reference pattern configuration satisfies certain conditions.

First of all, consider the reverse transition from E^{d+1} to E^d ; the $(d+1)^{\text{th}}$ feature is removed to reduce dimensionality from $d+1$ to d . The following assumptions on reference pattern configuration in E^{d+1} enable us to obtain a sufficient condition for the misclassification probability associated with the q^{th} orientation to decrease with the transition from E^{d+1} to E^d .

Let \underline{a}_{ij} be the ideal, or noise-free reference pattern associated with class ω_i , $i=1, 2$, at the j^{th} orientation. Let a total of s_1 discrete orientations of ω_1 have \underline{a}_{1q} as the ideal reference pattern. Similarly, let a total of s_2 discrete orientations of ω_2 have \underline{a}_{2q} as the ideal reference pattern. Noisy patterns corresponding to \underline{a}_{ij} are assumed

to be distributed according to $n(\underline{a}_{ij}, \sigma^2 I)$. One further assumption on reference pattern configuration in E^{d+1} is made so as to simplify the analysis. Let there be no other reference patterns of either class within the two $(d+1)$ -dimensional hyperspheres of radius $r=3\sigma$ centered at \underline{a}_{1q} and \underline{a}_{2q} . Fig. 4-7 illustrates a reference pattern configuration in E^2 where \underline{a}_{1q} and \underline{a}_{2q} satisfy the necessary conditions.

The decision surface between ω_1 and ω_2 in the vicinity of the q^{th} subclass is the hyperplane defined by

$$s_1 \cdot \exp \left[-\|\underline{x} - \underline{a}_{1q}\|^2 / 2\sigma^2 \right] = s_2 \cdot \exp \left[-\|\underline{x} - \underline{a}_{2q}\|^2 / 2\sigma^2 \right],$$

that is,

$$(4-7) \quad \underline{W}^t \underline{x} + W_0 = 0$$

where $\underline{W} = (1/\sigma^2) \cdot (\underline{a}_{1q} - \underline{a}_{2q})$ and $W_0 = (1/2\sigma^2) \cdot \|\underline{a}_{2q} - \underline{a}_{1q}\|^2 + \log(s_1/s_2)$.

The decision surface that minimizes the misclassification probability between the q^{th} subclass of each class is the hyperplane that is perpendicular to, and at the midpoint of, the line joining \underline{a}_{1q} and \underline{a}_{2q} . Now, if $(s_1/s_2)=1$ then Eq. (4-7) defines the decision surface that minimizes the misclassification probability between the q^{th} subclass of each class. In other words Eq. (4-7) with $s_1=s_2$ defines the optimal decision surface for the q^{th} orientation.

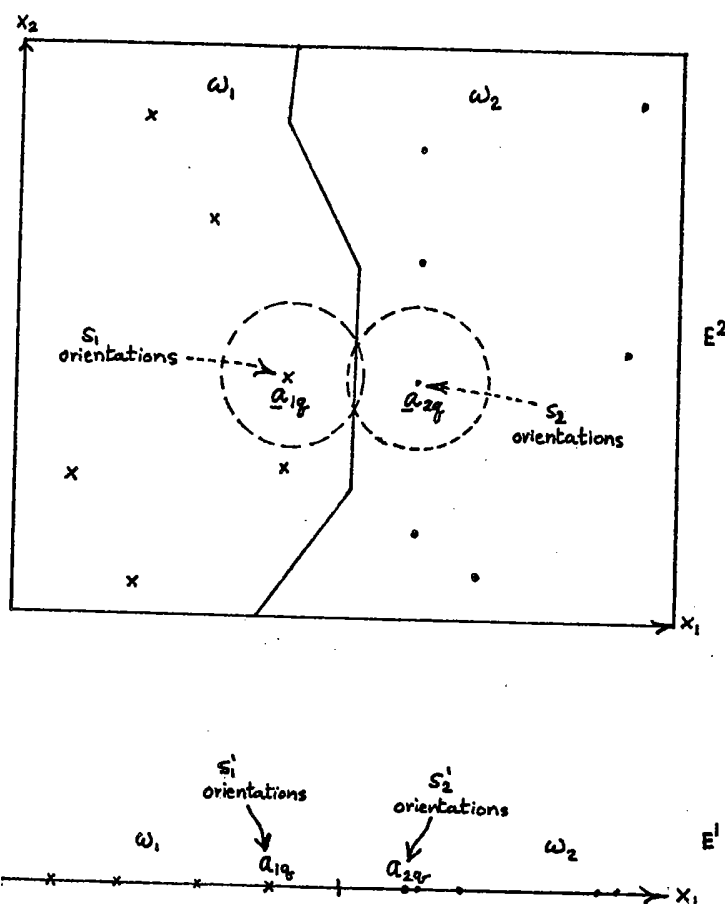


Fig. 4-7. Configuration of reference patterns to illustrate a negative effect of dimensionality.

If the transition from E^{d+1} to E^d leads to s_1' reference points to be projected onto \underline{a}_{1q} , and s_2' reference points onto \underline{a}_{2q} , with no other reference points falling within the two d -dimensional hyperspheres of radius r centered at \underline{a}_{1q} and \underline{a}_{2q} , then the decision surface between the q^{th} subclass of each class in E^d is the hyperplane defined by the above equation, except for W_0 , which is

given by

$$w_o = (1/2\sigma^2) \cdot \|\underline{a}_{2q} - \underline{a}_{1q}\|^2 + \log(s'_1/s'_2) \quad .$$

If the ratio (s'_1/s'_2) is closer to unity than the ratio (s_1/s_2) , it implies that the resulting decision surface is closer to the optimal decision surface for the q^{th} subclass. That is, the transition from E^{d+1} from E^d has changed the Sebestyen decision surface in such a way that as far as the q^{th} orientation is concerned the decision surface in E^d is closer to the optimal decision surface than the decision surface in E^{d+1} . In other words, the misclassification probability at the q^{th} orientation, or between the q^{th} subclass of each class, can be expected to decrease with the transition from E^{d+1} to E^d . It follows that the misclassification probability can increase with the transition from E^d to E^{d+1} .

The foregoing analysis concerning misclassification probability allows us to state a sufficient condition for the negative effect of dimensionality to occur. If α_{iq} is the value of the reference pattern of the q^{th} subclass of ω_i at the $(d+1)^{\text{th}}$ feature, say x_{d+1} , the following conditions on the feature-space configuration, and discriminating property of x_{d+1} will lead to an increase in

the misclassification property associated with the q^{th} subclass:

- (i) \underline{a}_{iq} is part of a cluster C_i in E^d , such that each of the C_i have approximately the same number of points, s_d ; the average intra-class separation of C_i , say w_i , is such that $w_i \ll \sigma^2$; and no reference patterns lie within distance 3σ from the centers of the clusters C_1 and C_2 ;
- (ii) feature x_{d+1} discriminates poorly at the q^{th} orientation as compared to its discriminating capacity at other orientations, i.e.,

$$|\alpha_{1q} - \alpha_{2q}| \ll |\alpha_{1j} - \alpha_{2j}|, \quad j=1, \dots, s, j \neq q;$$
- (iii) $\alpha_{ik} \cong \alpha_{iq}$ for either $i=1$ or $i=2$ and for s_{d+1} values of k , where $s_d \geq s_{d+1} \geq 1$.

The first condition assures that the decision surface is close to the optimal decision surface for the q^{th} orientation in E^d . Due to the second condition, the relative positions of the reference patterns of the q^{th} orientation remain relatively close to each other with the change in dimensionality. Due to the third condition a disproportionate number of points of each of the clusters move away from their original relative positions, which

leads to the Sebestyen classifier decision boundary to move toward the cluster with fewer reference points in E^{d+1} .

If feature x_{d+1} is introduced while condition (i) does not hold, i.e., the \underline{a}_{iq} are in arbitrary positions, we can expect a marginal decrease in the associated misclassification probability as a result of the small increase in the separation of the \underline{a}_{iq} . If the weight of a cluster is the ratio of the number of reference patterns in the cluster to the total number of reference patterns, then we can say the following concerning the dimensionality of the Sebestyen classifier: if the discriminating capacity of the feature to be introduced is small for the orientation possibility of interest, τ_q , and is large for other orientation possibilities, then a deterioration in performance can be expected if the original feature-space is such that the reference patterns \underline{a}_{iq} corresponding to τ_q are part of equally, and heavily, weighted clusters C_i .

IV.E.2. Dimensionality and the Nearest-Neighbor Classifier

The misclassification probability associated with the q^{th} subclass of each class can increase with dimensionality in the case of the nearest-neighbor classifier also. If p_q

is the average misclassification probability associated with the q^{th} subclass, the averaging being done over all classes, then p_q can be written in terms of its component probabilities as

$$p_q = \sum_{i=1}^S P(\text{not choose } \omega_i | \underline{a}_{iq} \text{ is true}) .$$

For p_q to increase with dimensionality, at least one of the component probabilities $P(\text{not choose } \omega_i | \underline{a}_{iq} \text{ is true})$ has to increase. The configuration and movement of reference patterns which leads to the increase of component probabilities will be examined here.

In the case of both the Sebestyen and nearest-neighbor classifiers, the effective probabilities of classifying a noisy version of \underline{a}_{iq} as not being from ω_i changes with dimensionality only when the following change occurs. A reference pattern lying within a hypersphere of radius r , and centered at \underline{a}_{iq} , changes its relative position with respect to \underline{a}_{iq} . The scalar r is effectively equal to 3σ for noise that is additive Gaussian, $n(0, \sigma^2)$, along each dimension.

The error probability conditioned on \underline{a}_{iq} will not increase with the introduction of a new feature as long

as not more than one reference pattern, other than \underline{a}_{iq} , lies within the above hypersphere in the original feature-space. Also, the error probability will not increase as long as all of the reference patterns within the hypersphere are classified as ω_i . The magnitude and sign of the change in $P(\text{not choose } \omega_i | \underline{a}_{iq})$ depends upon the relative locations of the reference patterns and upon the movement of the reference patterns with the introduction of the new feature.

Consider the univariate, two-class problem where $a_1 \in \omega_1$ is the subclass mean under consideration. Let $P(\omega_2 | a_1)$ represent the probability of the nearest-neighbor classifier choosing class ω_2 when a_1 is the true subclass mean. Let $b_1 \in \omega_2$ be a second distinct reference pattern that lies within the interval $(a_1 - r, a_1 + r)$. We will consider a third distinct reference pattern within this interval, and consider different possibilities for its location with respect to a_1, b_1 . The third reference pattern will either be $b_2 \in \omega_2$, or $a_2 \in \omega_1$, and occur in one of three possible positions with respect to a_1 and b_1 . The possible locations, assuming $b_1 > a_1$, are: to the left of a_1 , between a_1 and b_1 , and to the right of b_1 . Let the six possible configurations of reference patterns be labelled as follows:

- (1) $a_2 a_1 b_1$, (2) $a_1 a_2 b_1$, (3) $a_1 b_1 a_1$, (4) $b_2 a_1 b_1$,
 (5) $a_1 b_2 b_1$, (6) $a_1 b_1 b_2$.

We will examine the effect of increasing dimensionality on $P(\omega_2 | \underline{a}_1)$ of the nearest-neighbor classifier for each configuration.

Without loss of generality, assume the subclass mean of interest, \underline{a}_1 , to be the origin of E^d . Let the coordinates of b_1 in E^2 be (b_1, b'_1) , and let K be the displacement of the third reference pattern along the new feature. If we assume $b'_1 = 0$, then the behavior of $P(\omega_2 | \underline{a}_1)$ with increasing K is as follows: with configuration 1, the probability decreases first and then increases monotonically; with configurations 2 and 3, it increases monotonically; with configurations 4 and 5, it decreases monotonically; and with configuration 6, it increases first and then decreases monotonically.

Fig. 4-8a illustrates the behavior of $P(\omega_2 | \underline{a}_1)$ with K for each configuration. If $b'_1 \neq 0$, $P(\omega_2 | \underline{a}_1)$ behaves similarly with K' , the distance of the third reference pattern from the line joining the origin to the point (b_1, b'_1) . The effect of non-zero b'_1 is to decrease $P(\omega_2 | \underline{a}_1)$ with the change in dimensionality.

In order to express $P(\omega_2 | \underline{a}_1)$ as a function of K , consider configuration 2, and the corresponding nearest-

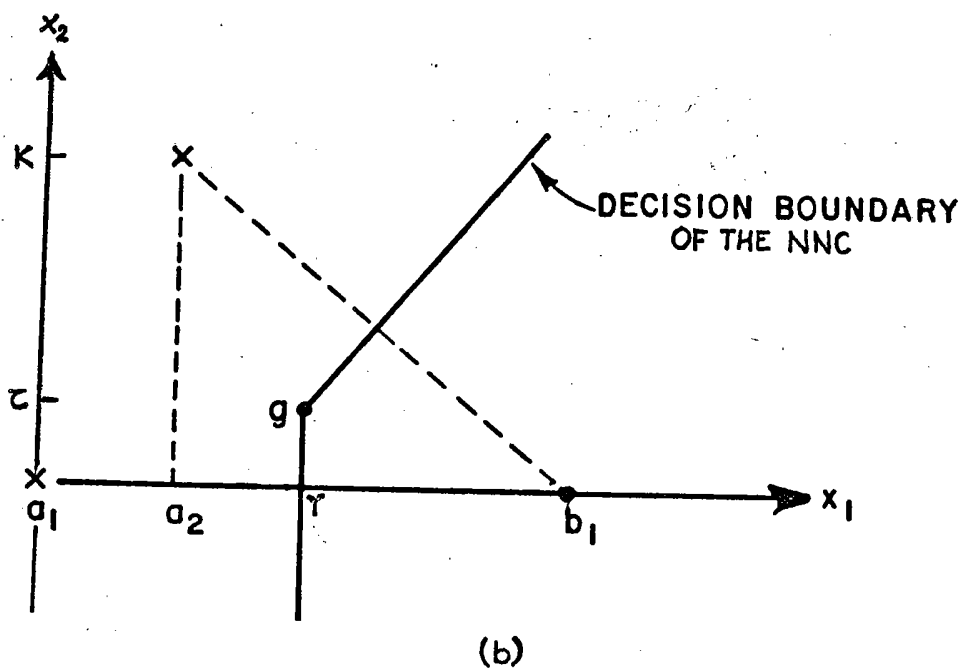
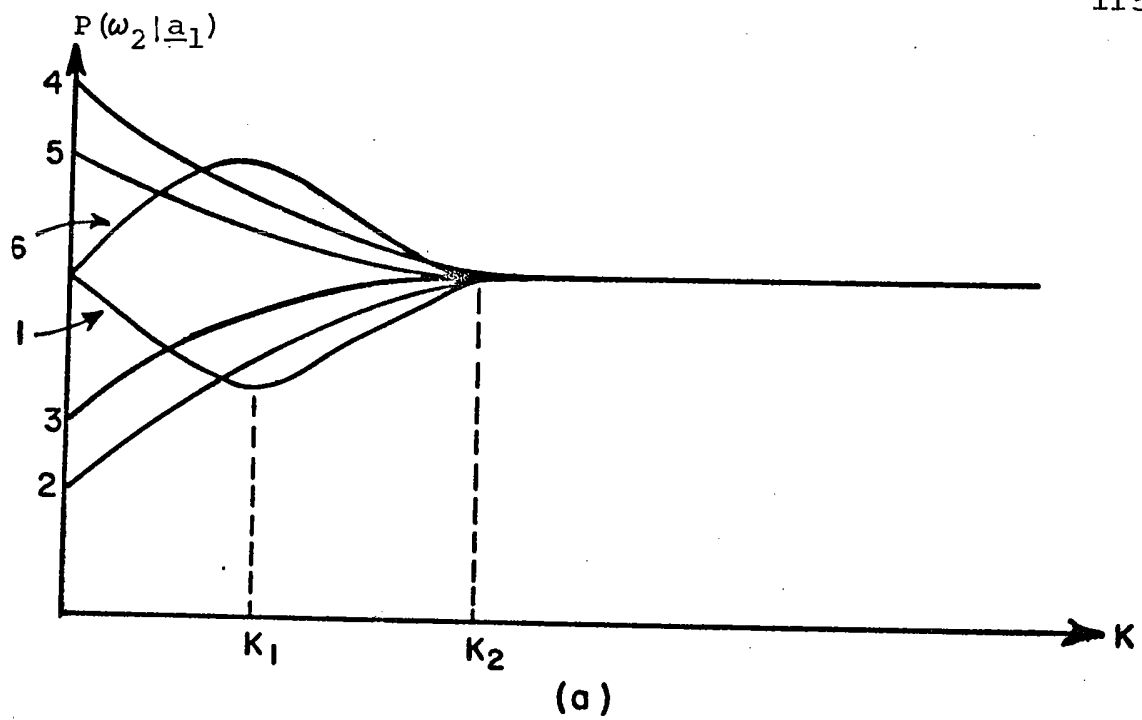


Fig. 4-8 Effect of dimensionality on the nearest-neighbor classifier.

neighbor decision boundary (Fig. 4-8b). Let g be the point of intersection of the pair of straight lines forming the decision boundary. Since g is also the center of the circle passing through the three reference points with coordinates $(0,0)$, (a_2, K) , and $(b_1, 0)$, the coordinates of $g \equiv (\tau, \epsilon)$ can be readily expressed as

$$\tau = (K/2) - (1/2K)(b_1 - a_2)a_2,$$

$$\epsilon = (b_1/2) .$$

At $K=0$, the decision boundary consists of a single straight line that is perpendicular to the real line at $(a_2 + b_1)/2$.

If $p(\underline{x}/\underline{a}_1)$ is the probability density function of the subclass associated with \underline{a}_1 , then

$$P(\omega_2/\underline{a}_1) = \int_{x_1=\alpha}^{\infty} \int_{x_2=-\infty}^{mx_1+\lambda} p(\underline{x}/\underline{a}_1) \cdot d\underline{x}$$

where m and λ , the slope and intercept of the skewed part of the decision boundary, are given by

$$m = (b_1 - \tau)/\epsilon, \quad \lambda = (\tau^2 + \epsilon^2 - \tau b_1)/\epsilon .$$

The expression $mx_1 + \lambda$ can be verified to be monotone increasing with K . That is, $P(\omega_2|\underline{a}_1)$ can be expected to be monotone increasing with K . In particular, if $K=0$, $p(\underline{x}/\underline{a}_1) \sim n(\underline{a}_1, \sigma^2 I)$, and $p(x_1/a_1)$, $p(x_2/a_2)$ are independent,

$$\begin{aligned}
 P(\omega_2/a_1) &= \int_{(a_2+b_1)/2}^{\infty} p(x_1/a_1) dx_1 \\
 &= (1/\sqrt{2\pi}) \left\{ \frac{\sigma}{(a_2+b_1)} \right\} \cdot \exp \left[-(a_2+b_1)^2 / 2\sigma^2 \right]
 \end{aligned}$$

which is the minimum value of the error probability for the transition from E^1 to E^2 . The maximum value occurs at $K=\infty$, for which

$$P(\omega_2/a_1) = \int_{b_1/2}^{\infty} p(x_1/a_1) dx_1$$

or,

$$P(\omega_2/a_1) \cong (1/\sqrt{2\pi}) \cdot (\sigma/b_1) \cdot \exp \left[-b_1^2 / 2\sigma^2 \right].$$

The error probability reaches 99% of its asymptotic value at

$$K = K_2 = (r^2 - b_1^2/4) \pm \left[(4r^2 - b_1^2) - (a_2^2 - a_2b_1)^{1/2} \right].$$

Similarly, for each of the other configurations, $P(\omega_2/a_1)$ can be expressed as a function of K . Thus the behavior of $P(\omega_2/a_1)$, as shown in Fig. 4-8a can be explained. In the case of configuration 6, the point of intersection of the straight lines forming the decision boundary of the nearest-neighbor classifier has a coordinate that decreases with increasing K , up to $K=K_1$, and then increases. The minimum value of τ occurs at

$$K = K_1 = \pm \sqrt{b_2 \cdot (b_2 - b_1)}.$$

So also, in the case of configuration 1, the maximum value of τ occurs at $K = K_1 = \pm \sqrt{a_2(a_2 - a_1)}$

We have thus been able to characterize, in an elementary situation, the change in misclassification probability of subclasses with the movement of reference patterns when the dimensionality of the nearest-neighbor classifier is increased. It has been observed that the misclassification probability of subclasses can increase with increased dimensionality. Thus if the nearest-neighbor classifier is employed in a problem where classes are characterized by reference patterns corresponding to s discrete orientations, the average misclassification probability associated with proper subsets of the s orientations can increase with dimensionality. This will happen if on the average there is an increase in the misclassification probabilities associated with orientations in the subset. It has also been observed that the misclassification probability associated with subclasses can increase with dimensionality even if the reference patterns are linearly separable in both the original and updated feature-spaces; examples are configurations 1, 2, and 6 in the foregoing analysis.

The utility of the analysis of dimensionality is primarily to serve as a warning to the user that misclassification probabilities of interest can actually increase with dimensionality, even when the reference patterns are noise-free. The linear separability of classes does not guarantee a decrease in the misclassification probability associated with particular orientations, as dimensionality is increased. As a consequence, a feature should be introduced only if the reference pattern configurations in the original and updated feature-spaces definitely indicate that the relevant misclassification probabilities will decrease.

CHAPTER V

CLASSIFICATION BIAS, MISCLASSIFICATION PROBABILITY, AND ECONOMY OF IMPLEMENTATION

The performance of competing classifier designs for radar aircraft identification are evaluated in this chapter. In any practical application of classification procedures one chooses performance measurement criteria that are appropriate for the application. The performance criteria that are of significance in implementing a pattern recognition system for radar aircraft identification are misclassification measures, economy of implementation, robustness with parameter variations, and extensibility to the multiple observation problem. In this chapter a comparison of misclassification probabilities of the Sebestyen and nearest-neighbor classifiers is made. The comparison is done experimentally using radar aircraft data. The economy of implementation of different classifiers are compared from the point of view of number of the computation steps required

to make decisions. The criterion of extensibility to the multiple observation problem will be discussed in Chapter VI.

V.A. Classification Bias

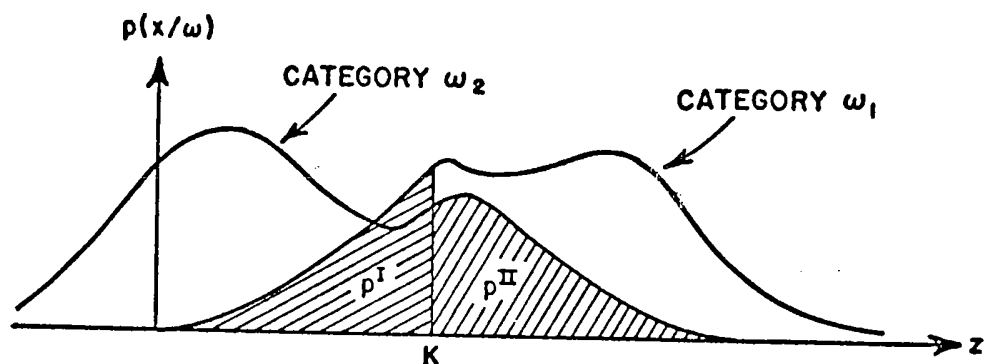
In order to compare the performance of two or more classifiers which have nearly identical average misclassification probabilities, we propose a probabilistic criterion here. The following discussion leads to the definition of the criterion, to be called classification bias, and its significance to classifier performance.

In a two-category pattern recognition problem, with categories ω_1 and ω_2 , we can define two error probabilities as follows:

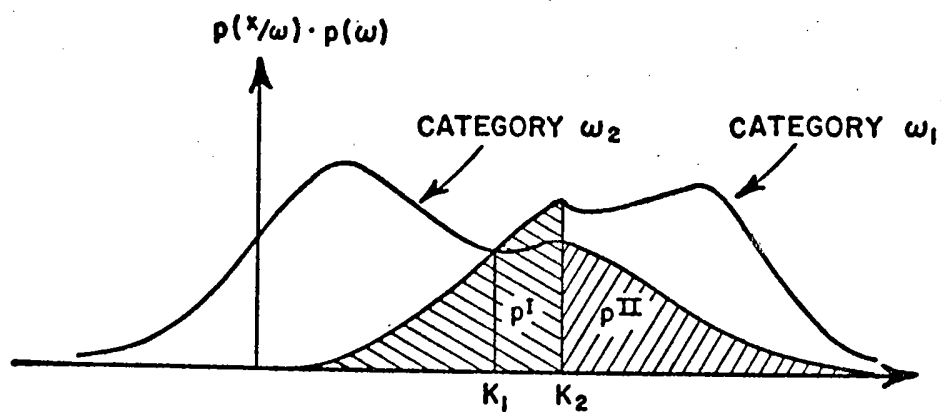
p_1 = probability of choosing ω_2 when the true state of nature is ω_1

p_2 = probability of choosing ω_1 when the true state of nature is ω_2 .

In the univariate case we can illustrate these probabilities in terms of class-conditional probability density functions as shown in Fig. 5-1a. The two error probabilities are known by different terminology in each of the literatures of statistical hypothesis testing and signal detection theory. In the former they correspond to type 1 and type 2 error probabilities if ω_1 corresponds to the null hypothesis



(a)



(b)

Fig. 5-1. Example of non-zero classification bias of the Bayes classifier. K_1 =Bayes decision boundary, K_2 =Decision boundary to minimize classification bias.

and ω_2 to the alternate hypothesis. Further, for a given hypothesis test, p_1 is referred to as the size of the test and $(1-p_2)$ as the power of the test. Similarly, in signal detection theory the two error probabilities are known as the probabilities of "false rest" and "false alarm" where category ω_1 corresponds to signal being present and category ω_2 to signal being absent. If $P(\omega_1)$ and $P(\omega_2)$ are the a priori probabilities of the two categories, then the average probability of misclassification of the classifier is:

$$(5-1) \quad P_e = P(\omega_1) \cdot p_1 + P(\omega_2) \cdot p_2 .$$

In designing the Neyman-Pearson hypothesis test, one fixes the value of p_1 and designs the test so as to minimize p_2 . Whereas, in designing a pattern classifier, one usually takes the Bayesian approach of minimizing P_e . However, neither the Neyman-Pearson test nor the Bayes classifier guarantees that the resulting p_1 and p_2 will be equal. Fig. 5-1b illustrates a situation where the optimal Bayes classifier does not result in $p_1 = p_2$.

Definition:

When a two-category classifier has associated class-conditional error probabilities p_1 and p_2 , with a priori

probabilities of class membership being $P(\omega_1)$ and $P(\omega_2)$, the classifier is said to have classification bias, B , given by

$$(5-2) \quad B = |P(\omega_1) \cdot p_1 - P(\omega_2) \cdot p_2| .$$

In other words, classification bias of a classifier is the weighted difference between its type 1 and type 2 error probabilities. It follows that for any classifier the following inequality is true.

$$0 \leq B \leq P_e \leq \frac{1}{2} .$$

B is equal to its upper bound P_e , if and only if either class a priori probability or either class-conditional error probability vanishes.

Consider now a decision problem where it is required to select one of two classifiers C_a and C_b . Each of C_a and C_b is designed for the same two-category classification problem. It is assumed that the average probability of misclassification of each of the two classifiers is the same. That is, if p_1^a, p_2^a are the class-conditional error probabilities of C_a , and p_1^b and p_2^b are the class-conditional error probabilities of C_b , then

$$\begin{aligned} P_e &= P(\omega_1) \cdot p_1^a + P(\omega_2) \cdot p_2^a \\ &= P(\omega_1) \cdot p_1^b + P(\omega_2) \cdot p_2^b . \end{aligned}$$

Further, it is assumed that $B_a < B_b$, where

$$B_a = |P(\omega_1) \cdot p_1^a - P(\omega_2) \cdot p_2^a|$$

$$B_b = |P(\omega_1) \cdot p_1^b - P(\omega_2) \cdot p_2^b|.$$

We suggest that factors such as computational complexity of classifiers being either equal or inconsequential, one ought to choose the less biased classifier, i.e., C_a . The reason for this choice is one of making the classifier symmetric with respect to the two classes. The first advantage of symmetry is simplicity. For example, when $P(\omega_1) = P(\omega_2)$ the average misclassification probability is itself the class-conditional misclassification probability. A second and more important advantage for radar aircraft identification occurs in the multiple decision situation which is described here briefly and dealt with in detail in Chapter VI.

Let C be a classifier designed to decide between two classes ω_1 and ω_2 . Let $P(\omega_1) = P(\omega_2) = (\frac{1}{2})$. Let p_1 and p_2 be the class-conditional error probabilities of C . Let x_1, \dots, x_n be n noisy patterns from the same, but unknown, class. Consider the following procedure: submit each of the n patterns to C and decide the class on the basis of a majority-vote over n decisions of C . If n is odd, then the

class-conditional error probability of this majority-vote procedure can be shown to be

$$p_i(n) = \sum_{j=\frac{n+1}{2}}^n \binom{n}{j} p_i^j (1-p_i)^{n-j}$$

and the average misclassification probability to be

$$P_e(n) = \frac{1}{2} \sum_{i=1}^2 p_i(n) .$$

The rate of decrease of $P_e(n)$ with n is dependent on the quantity $|p_1 - p_2|$. For a constant value of $P_e = (1/2) \cdot (p_1 + p_2)$ and any given n , $P_e(n)$ takes on increasing values for increasing values of $|p_1 - p_2|$. A proof of this result is given in Chapter VI. As a consequence, the number of patterns required to achieve a specified error probability increases with classifier bias.

V.B. Estimating the Misclassification Probability of the Nearest-Neighbor Classifier

From a theoretical standpoint, the misclassification probability of each of the nearest-neighbor classifiers discussed in Section IV.B can be analytically determined for a given set of reference patterns, and specified $p(\underline{x}/\omega_i)$ and $P(\omega_i)$. However error probability calculations are difficult, although the concept is quite simple. If

p_i is the probability that a classifier does not choose class ω_i when ω_i is the true state of nature, i.e., p_i is the class-conditional probability of error, and $P(\omega_i)$ is the class a priori probability, then the average misclassification probability of the classifier is given by

$$(5-3) \quad P_e = \sum_{i=1}^c P(\omega_i) \cdot p_i \quad .$$

The class-conditional error probability of a given classifier can be written as

$$(5-4) \quad p_i = 1 - \int_{R_i} p(\underline{x}/\omega_i) d\underline{x}$$

where R_i are mutually exclusive partitions of the feature-space such that the classifier chooses class ω_i when \underline{x} is in R_i . Substituting (5-4) into (5-3), the average misclassification probability can be written as

$$(5-5) \quad P_e = 1 - \sum_{i=1}^c \int_{R_i} p(\underline{x}/\omega_i) P(\omega_i) d\underline{x}$$

In the case of the nearest-neighbor classifier based on the Euclidean distance measure, region R_i associated with class ω_i is defined by

$$R_i = \left\{ \underline{x} \mid \min_{j=1, \dots, s} \left[\|\underline{a}_{ij}\|^2 - \frac{1}{2} \underline{a}_{ij}^t \underline{x} \right] < \min_{\substack{j,k=1, \dots, s \\ k=(i+1), \dots, c}} \left[\|\underline{a}_{kj}\|^2 - \frac{1}{2} \underline{a}_{kj}^t \underline{x} \right] \right\} \quad .$$

For arbitrary values of c and s , the boundaries of the regions R_i are too complex for the task of evaluating P_e by the multidimensional integration approach of Eq. (5-5). For example, to characterize the piece-wise linear decision boundary of the nearest-neighbor classifier based on the Euclidean distance, a total of cs hyperplanes may be required. Thus, in order to determine the misclassification probability of the nearest-neighbor classifier, we use an experimental technique.

Many techniques for the experimental evaluation of classifiers are discussed in the literature (Toussaint, 1974): Most of the techniques, such as the "hold-out" method, and "leave-one-out" method are specifically designed for the situation where only a finite number of sample patterns are available to design and test the classifier. However, in evaluating classification procedures for the radar aircraft identification problem, a potentially infinite number of test samples are available, since they can be generated from the reference patterns in accordance with the specified noise distributions $p(\underline{x}/\omega_i, J)$. The experimental evaluation technique is to generate noisy patterns according to the specified

distribution by using a random number generator. These noisy patterns are submitted to the classifier to be evaluated. It can be shown that the maximum-likelihood estimate of the misclassification probability of the classifier is given by the ratio of the number of incorrectly classified patterns to the total number of test patterns (see Highleyman, 1962). Such an estimate converges to the true misclassification probability as the number of samples tends to infinity.

In the experimentation, nine discrete orientations were assumed to be equally likely, two classes possible, and twenty experiments performed at every possible orientation. This corresponds to a total of 360 experiments for a specified noise standard deviation. The confidence associated with such an estimate can be determined from the curves provided by Highleyman (see Duda and Hart, 1973). For example, if by using 360 experiments the misclassification probability is estimated to be 0.3, then the true misclassification probability lies between 0.26 and 0.34 with 95% confidence. An attempt to decrease the variance of the estimate by half would mean that four times as many experiments would have to be conducted, i.e., 1440 experiments in

this case. Since there are other approximations in the overall aircraft identification model, 20 experiments per discrete orientation was regarded to be reasonable. A program to evaluate each of the classifiers discussed in Section III.B has been developed. For a description of the program, and a more comprehensive set of results see Srihari and White (1975).

V.C. Experimental Evaluation of the Nearest-Neighbor and Bayes Classifiers

The average misclassification probability and classification bias of the Sebestyen classifier, nearest-neighbor classifier, and a non-Euclidean nearest-neighbor classifier were evaluated with data from the radar aircraft identification problem. The misclassification characteristics as functions of noise standard deviation are given in Figs. 5-2 and 5-3. The characteristics correspond to the problem of discriminating between the F-104 and the F-4. The true orientation in this case was approximately $(0^\circ, 30^\circ)$, a typical orientation. The measured orientation (ϕ_o, θ_o) was taken to be $(0^\circ, 30^\circ)$ and $q(\phi, \theta / \phi_o, \theta_o)$ was regarded to be uniform over nine discrete orientations in the range $\phi \pm 15^\circ, \theta \pm 10^\circ$.

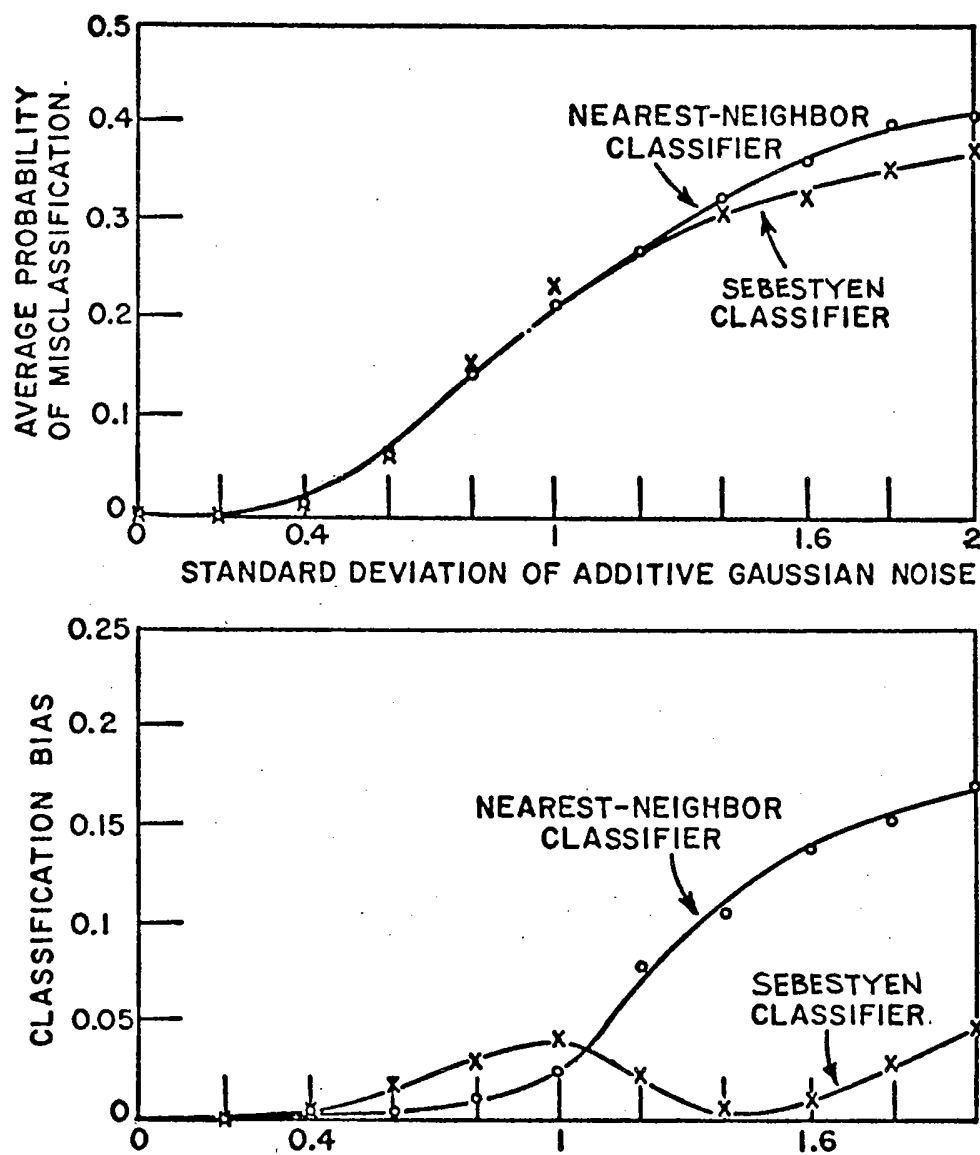


Fig. 5-2. Evaluation of classifier performance with radar aircraft data. Classes: F104, F4 pair; $s=9$, in the vicinity of $(0^\circ, 30^\circ)$; amplitude and phase data.

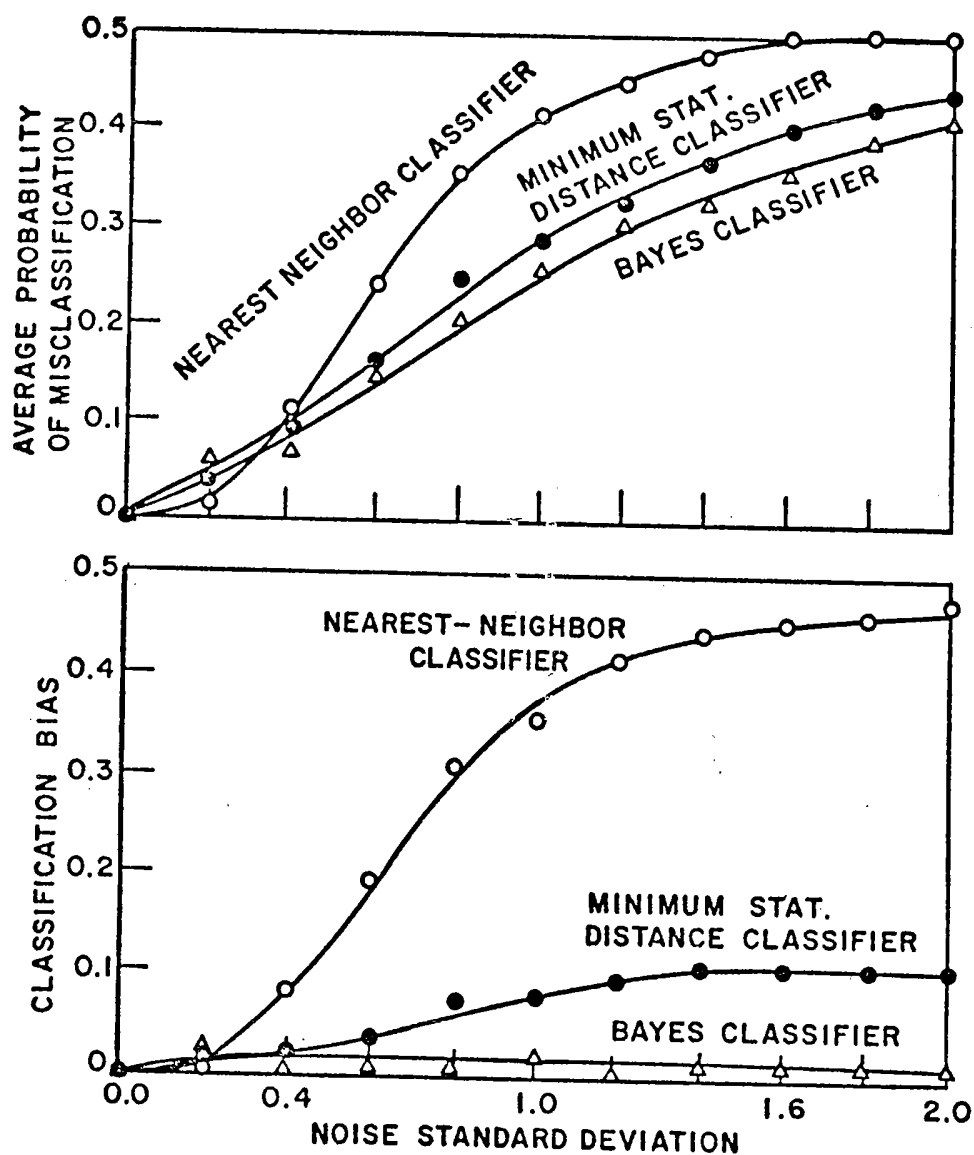


Fig. 5-3. Evaluation of classifier performance with radar aircraft data. Classes: F-104, F-4 pair; $s=9$, in the vicinity of $(0^\circ, 30^\circ)$; amplitude data only.

The characteristics of Fig. 5-2 correspond to the case where class-conditional densities are Gaussian mixtures with $s=9$ and $d=8$, i.e., the pattern vector was constructed using both amplitude and phase data at four different frequencies. The noise standard deviation used was a normalized value, the normalization factor being the average amplitude return over all four frequencies and both classes. When class-conditional densities are Gaussian mixtures, the minimum statistical-distance classifier is identical to the nearest-neighbor classifier, thus only two sets of characteristics are shown. The average misclassification probabilities of the Sebestyen classifier and the nearest-neighbor classifier tended to be nearly the same over all values of noise standard deviation considered (Fig. 5-2a). This behavior was found to be true in most of our experimentation with data from the radar aircraft identification problem (see Ref. 30). The explanation lies in the fact that reference patterns are situated such that the inter-class separation is considerably smaller than the intra-class separation. Note also that the misclassification probability of the Sebestyen classifier rises steeply with

noise standard deviation up to a value near unity, after which its rate of increase is small. That is, from Theorem 3.1, the minimum inter-class separation between the F-4 and F-104 in this case is 50% greater than the average amplitude return.

In spite of the fact that the average misclassification probabilities of the Sebestyen classifier and the nearest-neighbor classifier tend to be the same, the values of their classification bias are different (Fig. 5-2b). The classification bias of the Sebestyen classifier is apt to fluctuate at a low absolute value. However, the classification bias of the nearest-neighbor classifier increases monotonically with noise standard deviation. An intuitive explanation of the monotonic increase in the classification bias of the nearest-neighbor classifier is as follows. Consider the nearest-neighbor decision surface, as shown for a particular case in Fig. 4-1b. As the noise standard deviation is increased the class-conditional error probability associated with class ω_1 increases at a faster rate than the class-conditional error probability associated with class ω_2 . This may be seen by considering circles of

radius 2σ centered at discrete points on each of the data surfaces. Since class ω_1 is on the concave side of the decision surface, an increase in σ leads to the following effect. Each circle centered on the ω_1 data surface has an increased proportion of its area which lies on the ω_2 side of the decision surface. Similarly, each circle centered on the ω_2 data surface has a decreased proportion of its area which lies on the ω_1 side of the decision surface. As a consequence the classification bias of the nearest-neighbor classifier increased monotonically with noise standard deviation. On the other hand, the classification bias of the Sebestyen classifier does not increase monotonically with noise standard deviation since its decision boundary varies with σ so as to minimize the average misclassification probability. The non-zero classification bias of the Sebestyen classifier is due to the non-symmetric nature of the class-conditional probability distributions.

The characteristics of Fig. 5-3 correspond to the situation where class-conditional probability distributions are Rician mixtures, $s=9$, and $d=4$, i.e., the feature vector consists of amplitude data at four different frequencies. In this case both misclassification probability and classification bias of the nearest-neighbor classifier are

considerably higher than that of the Bayes classifier. The performance characteristics of the minimum statistical-distance classifier, which is C_{mult} with Rician $p(\underline{x}/\omega_{ij})$ (see Section IV.B), are close to that of the Bayes classifier. This is explained by the fact that the relationship between the two classifiers is similar to that between the Sebestyen and nearest-neighbor classifiers when class-conditional densities are Gaussian mixtures.

V.C.1. Robustness of the Sebestyen Classifier

The robustness of a classifier is a measure of how well the classifier performs when the class-conditional distribution parameters differ from the classifier assumption of parameters. In Section IV.D the robustness of the Sebestyen classifier was examined analytically in a restrictive univariate case. The examination was from the point of view of changes in the Sebestyen classifier decision boundary when the noise variance assumption is changed from its true value. In this section the robustness of the Sebestyen classifier with the assumed value of noise standard deviation is experimentally examined in the multivariate case. Since the decision boundary in the multivariate case is a

hypersurface, a convenient scalar measure of the effect of an incorrect assumption of noise variance is the average misclassification probability. That is, we will look at the changes in misclassification probability of the Sebestyen classifier when the assumed noise standard deviation, σ_a , differs from the true value σ_t .

In terms of σ_a and σ_t , the average misclassification probability of the Sebestyen classifier in the two-class case is expressed as

$$P_e = 1 - \frac{1}{(2\pi)^d \sigma_t^{2d}} \sum_{i=1}^2 \sum_{k=1}^s \int_{R_i} \exp\left[-\|\underline{x} - \underline{a}_{ik}\|^2 / (2\sigma_t^2)\right] \cdot P(\omega_i) d\underline{x}$$

where the mutually exclusive regions R_1 and R_2 of the feature-space are separated by the hypersurface defined by

$$\sum_{k=1}^s \exp\left[-\|\underline{x} - \underline{a}_{1k}\|^2 / (2\sigma_a^2)\right] - \sum_{k=1}^s \exp\left[-\|\underline{x} - \underline{a}_{2k}\|^2 / (2\sigma_a^2)\right] = 0.$$

Experimental estimates of P_e are given in Fig. 5-4, where SIG refers to the assumed value of noise standard deviation σ_a , $s=9$ in the vicinity of $(0^\circ, 30^\circ)$, and $d=8$. It can be observed that $\sigma_a = \sigma_t - \delta$ leads to fewer errors than $\sigma_a = \sigma_t + \delta$ for σ_a values within the region of interest, viz., $(0, 2)$. This experimentally supports the theoretical developments for the univariate case in Section IV.D. A comparison of the characteristics of Figs. 5-2 and 5-4 shows that the

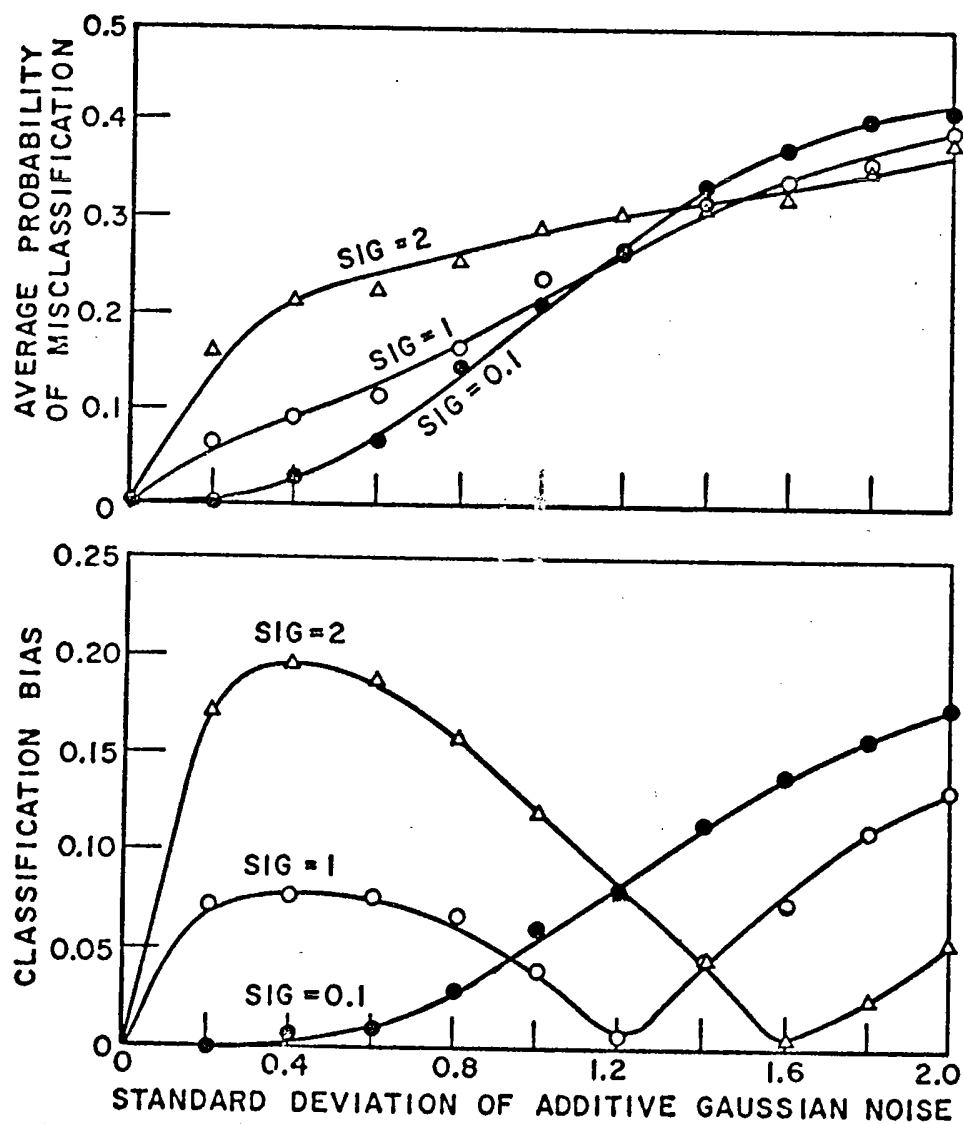


Fig. 5-4. Robustness of the Sebestyen classifier. Abscissa is the true noise level, SIG is the assumed noise level.

Sebestyen classifier with an assumption of $SIG=0.5$ approximates the performance of the nearest-neighbor classifier. This value corresponds to half the average amplitude return from both classes over all frequencies.

Misclassification estimates of the Sebestyen classifier when class-conditional densities are Rician mixtures are given in Fig. 5-5. The characteristics were obtained when the Sebestyen classifier assumption of noise variance is the same as the noise variance used to generate the Rician random variable. However, note that the Sebestyen classifier is not the Bayes classifier for the problem. Here again, the performance of the Sebestyen classifier is similar to that of the nearest-neighbor classifier. This can be seen by comparing the nearest-neighbor classifier characteristics with the Sebestyen classifier characteristics. Based on the empirical evidence, we expect the performance of the Sebestyen classifier to be similar to the nearest-neighbor performance, when subclass distributions are Rician and a reasonable assumption for σ is made.

V.C.2. Effect of Increasing the Number of Subclasses

The effect of increasing orientation uncertainty, or

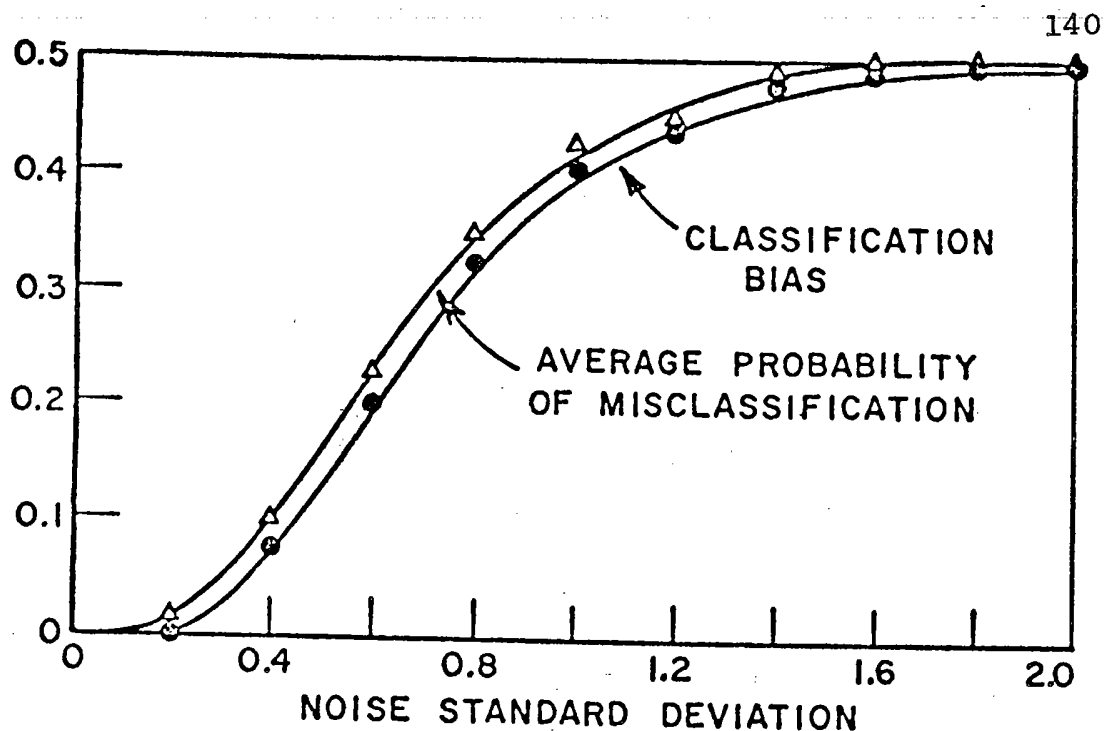
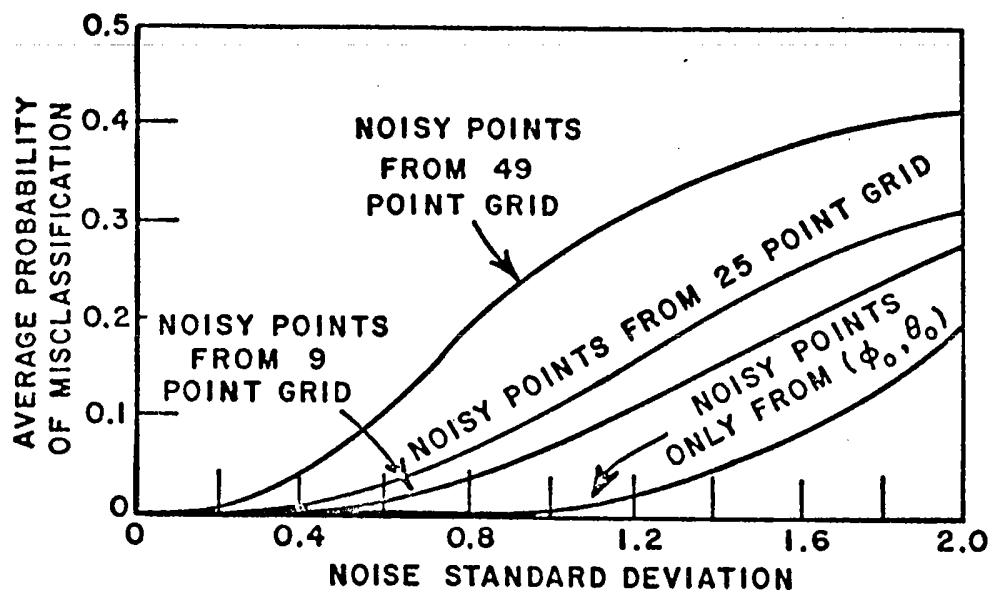
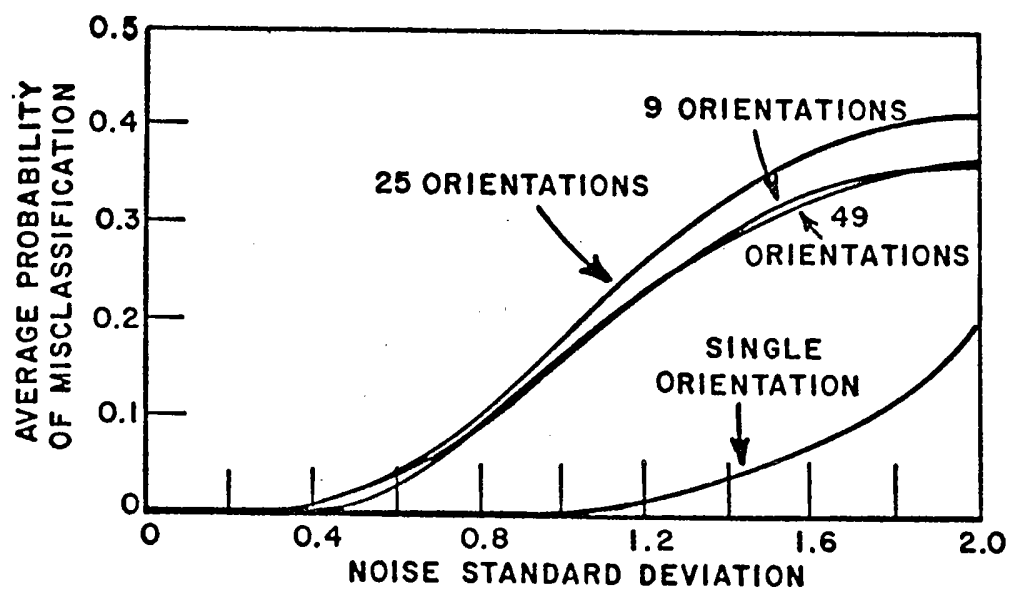


Fig. 5-5. Performance of the Sebestyen classifier when true class-conditional densities are Rician mixtures.

the number of possible subclasses on misclassification is examined here. The effect may be regarded as one of decreasing knowledge concerning the orientation of the aircraft. The effect is experimentally examined here from two viewpoints. Assuming a region of the orientation-space, $R(\theta_o, \theta_o)$ centered at (θ_o, θ_o) , we examine the variation of average misclassification probability as the size of R is increased. It is assumed that the resolution of quantization of the orientation space is fixed. Fig. 5-6a shows the effect of increasing the number of subclasses on the overall



(a)



(b)

Fig. 5-6. Effect of increasing the number of subclasses on
 (a) overall misclassification
 (b) misclassification associated with a particular subclass.

average misclassification probability when using the nearest-neighbor classifier, Gaussian noise, the F-104: F-4 pair, and $(\phi_0, \theta_0) = (0^\circ, 90^\circ)$. Each of the curves corresponds to a different value of s , the number of subclasses. The monotone increase of the overall average misclassification probability, P_s^* , with s can be explained as follows. In Section III.E it was shown that

$$P_s^* \geq (1/s) \sum_{i=1}^s p_i^*$$

where p_i^* is the misclassification probability between the i^{th} subclass of each class. P_s^* can then be written as

$$(5-6) \quad P_s^* = (1/s) \cdot \sum_{i=1}^s p_i^* + \delta$$

where δ is a small positive quantity. Similarly we approximate P_{s+1}^* as

$$(5-7) \quad P_{s+1}^* = \left[1/(s+1)\right] \cdot \sum_{i=1}^{s+1} p_i^* + \delta$$

The approximation is justified by the fact that the difference between the overall misclassification probability and the average of the individual misclassification probabilities is small. Eq. (5-6) and (5-7) can be combined and rewritten as

$$(5-8) \quad P_{s+1}^* = P_s^* \left[\frac{s}{s+1} \right] + \left[P_{s+1}^* / (s+1) \right] + \delta \cdot \left[1 - \frac{s}{s+1} \right].$$

For large values of s , the ratio $s/(s+1)$ can be replaced by unity in Eq. (5-8), i.e.

$$(5-9) \quad P_{s+1}^* = P_s^* + [p_{s+1}^*/(s+1)] .$$

Since $p_{s+1}^*/(s+1)$ is a positive quantity, it follows that for large s , P_s^* is a non-decreasing function of s .

When the number of subclasses is increased, the misclassification probability associated with a particular subclass is likely to change. This change depends on the reference pattern configuration. The characteristics of Fig. 5-6b correspond to the average misclassification probability of the particular orientation $(0^\circ, 90^\circ)$ as the number of reference patterns in the vicinity of $(0^\circ, 90^\circ)$ is increased. As the number of orientation possibilities is increased from 25 to 49, the misclassification probability associated with $(0^\circ, 90^\circ)$ is seen to decrease. An intuitive explanation of this phenomenon is given in Fig. 5-7. The decision surface of the classifier changes as the number of discrete orientations in $R(\theta_o, \theta_o)$ is increased. Assuming constant quantization of the orientation-space, the distribution defined by $n[f_i(\theta_o, \theta_o), \sigma^2 I]$ intersects with a varying decision surface. If the data surfaces are highly convoluted, the region of feature-space contributing

to the misclassification probability of the subclass under question can increase. Fig. 5-7a corresponds to the case where the data surfaces considered correspond to the region of orientation-space defined by $R_1(\varnothing_o, \theta_o)$. The decision boundary shown is that of the nearest-neighbor classifier. The region of feature-space that contributes effectively to the misclassification probability of the nearest-neighbor classifier, when the given class is ω_1 and the orientation $(\varnothing_o, \theta_o)$, is indicated as the shaded region of the figure. This region was obtained by the intersection of the effective non-zero probability contour with the nearest-neighbor decision surface. Fig. 5-7b corresponds to a set of orientations $R_2(\varnothing_o, \theta_o)$ which is a superset of $R_1(\varnothing_o, \theta_o)$. The data surfaces in this case are extended forms of those in Fig. 5-7a. Due to the convolutions of the data surfaces, the nearest-neighbor decision boundary changes in such a way that the region of misclassification associated with class ω_1 at $(\varnothing_o, \theta_o)$ is decreased. Consequently, decreased misclassification probability is observed for a particular class and orientation possibility, when orientation uncertainty is increased. If class ω_2 at $(\varnothing_o, \theta_o)$ is also subject to the same phenomenon, then the average

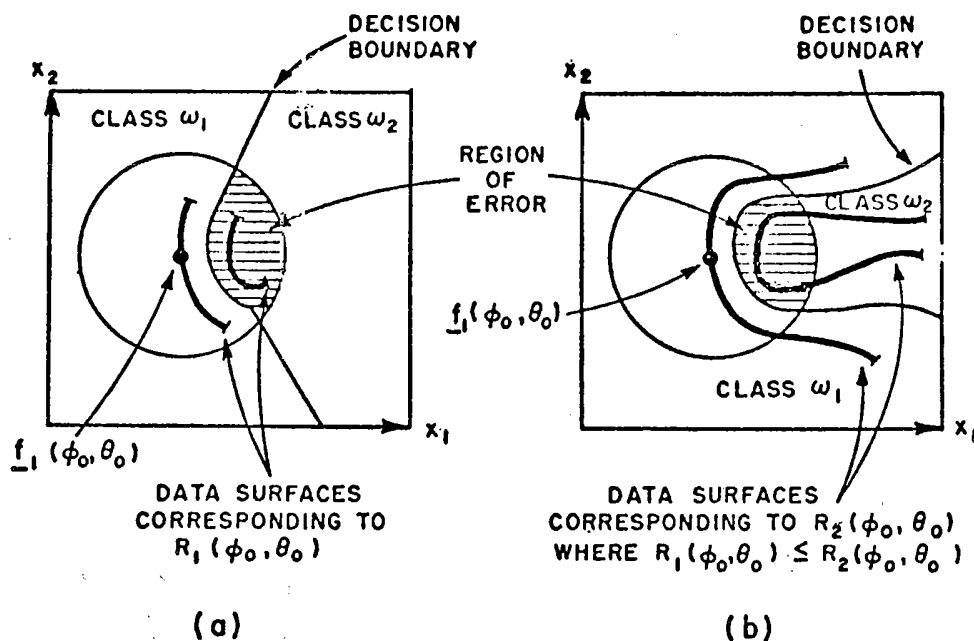


Fig. 5-7. Variation of nearest-neighbor decision surface with increased number of orientation possibilities.

misclassification probability at (ϕ_0, θ_0) will decrease when orientation uncertainty is increased.

V.D. Economy of Implementation

The complexity of calculation is of importance in any real-time application of a pattern recognition algorithm, especially in applications such as radar aircraft identification. In this section we examine the computational complexity of each of the classification algorithms discussed in earlier sections.

Since the quantities V_{ij} , $d_{ij}(\underline{x})$ and $d'_{ij}(\underline{x})$ are common to more than one algorithm (see Section IV.B), we consider first the complexity of computing these quantities.

TABLE 1

Quantity	Multiplications	Additions	Logarithms	Calls to Bessel Function Routines
V_{ij}	d	$d-1$	0	0
$W_{ij}^t \cdot \underline{x}$	d	$d-1$	0	0
$d_{ij}(\underline{x})$	$2d$	$3d-2$	0	0
$d'_{ij}(\underline{x})$	$2d$	$3d-2$	d	d

Assuming c classes and s discrete orientations, the computational complexity of each of the algorithms can be summarized as in Table 2.

Since the storage requirements of each of the five classifiers are the same, each of them requiring the storage of cs vectors, we can compare the computation time required in making a decision using each one of the classifiers. The relative time requirements depend upon the time taken to perform the elementary operations of addition, multiplication

TABLE 2

Classifier	Multiplications	Additions	Logarithms	Bessel Calls	Exponents	Comparisons
Nearest-neighbor	2dcs	(3d-2)cs	0	0	0	cs-1
Weighted-nearest-neighbor	2dcs	3dcs	0	0	0	cs-1
Minimum stat. distance (Rician)	2dcs	3dcs	dcs	dcs	0	cs-1
Sebestyen	(2d+1)cs	3dcs	0	0	cs	c-1
Bayes (Rician subclasses)	(2d+1)cs	3dcs	dcs	dcs	cs	c-1

and comparison. The relative time requirements for the elementary operations depend upon the architecture of the computer on which they are implemented. We compare here the relative time requirements of the nearest-neighbor and Sebestyen classifiers when implemented on a typical general purpose computer such as the DEC PDP-10.

Using fixed-point arithmetic, the multiplication operation on the PDP-10 requires 10.6 microseconds which is approximately four times the time taken for either an addition or a comparison operation (2.5 microseconds). Assuming the exponentiation operation to require an average of ten multiplications, the relative time requirements of the nearest-neighbor and Sebestyen classifiers in making a single decision can be expressed as follows:

$$\begin{aligned} T_{nnc} &= 4(2dcs) + (3d-2)cs - 1 \\ &= cs(11d-1) - 1 \end{aligned}$$

$$\begin{aligned} T_{sc} &= 4(2d+1)cs + 3dcs + 40cs + (c-1) \\ &= cs(11d+44) + c - 1 \end{aligned}$$

where the add operation is assumed to take unit time. The time requirement for the Sebestyen classifier is in the worst-case situation where all exponentials are evaluated to determine the discriminant function values. If the Sebestyen classifier were implemented using Algorithm B

that was developed in Section III.C, then the Sebestyen classifier can require as few as cs comparisons instead of cs exponentiations and $(c-1)$ comparisons.

The relative efficiency of implementation of the two classifiers is expressed by the ratio (T_{sc}/T_{nnc}) . Substituting the above expressions for T_{sc} and T_{nnc} into the ratio, we can write

$$T_{sc}/T_{nnc} = 1 + (4/d) \quad .$$

The number of classes and number of reference patterns have negligible effect on the ratio; their effect is felt most strongly when c and s are small. That is, the relative execution times of the Sebestyen and nearest-neighbor classifiers are determined primarily by the dimensionality of the feature-space. At worst, the Sebestyen classifier takes five times the execution time of the nearest-neighbor classifier, which happens when only a single feature is used. With eight features, as was used in our experimentation for the radar aircraft identification problem, the Sebestyen classifier takes 50% more execution time than the nearest-neighbor classifier. We can thus say that for large values of d , there is little to be gained computationally by using the nearest-neighbor classifier over the Sebestyen classifier.

CHAPTER VI

CLASSIFICATION BASED ON MULTIPLE OBSERVATIONS

If classification can be based on the cumulative evidence collected from n repeated observations, the probability of misclassification can be reduced from that obtained with a single observation. In this chapter we will examine alternative methods of utilizing the nearest-neighbor and Bayes classifiers for the multiple observation problem. A given classifier can be utilized in the multiple observation situation in a number of suboptimal ways. Such procedures are suboptimal in the sense that the minimum misclassification probability for the problem may not always be achieved. We will consider one such suboptimal procedure in detail. The procedure is one that chooses the class by taking a majority-vote over n different observations of a given classifier. In other

words, the procedure submits each of n observations to a given classifier and chooses that class which is represented most often among n decisions. The performance of the majority-vote procedure will be evaluated for different classifiers.

The optimal procedure for the problem, with constant orientation and Gaussian subclass distribution, is conceptually simple. This procedure, which consists of submitting the sample mean to the Bayes classifier, will be considered first. It will be followed by a description of the suboptimal procedure, its error analysis, and the comparison of the two procedures with radar aircraft data. Error analysis of the procedure consists of characterizing the behavior of misclassification probability with parameters such as number of samples, and class-conditional error probabilities of the classifier. In the second part of the chapter we will consider the problem of sequential classifier design for the problem, as opposed to fixed sample-size procedures considered in the first part of

the chapter. We will show how the Sebestyen classifier can be utilized in designing a sequential classifier, and finally evaluate the performance of such a procedure with the radar aircraft identification problem.

VI.A. An Estimate-Based Procedure

Consider the fixed sample-size hypothesis testing problem where $\underline{x}_1, \dots, \underline{x}_n$ are n observed patterns from an aircraft which are to be used in determining the classification of the aircraft. If aircraft orientation is assumed to remain constant, i.e., the samples are identically distributed, then the problem is one of composite hypothesis testing. The hypotheses correspond to location of noise-free patterns on different data surfaces.

If the observed patterns are stochastically independent, s discrete orientations are equally likely, and $p(\underline{x}/\omega_{jk}) \sim n(\underline{a}_{jk}, \sigma^2 I)$ is the distribution of the j^{th} class at the k^{th} orientation, then the problem optimal procedure can be shown to be equivalent to submitting the maximum-

likelihood estimate of the noise-free pattern, viz.,

$$(6-1) \quad \hat{\underline{x}} = \sum_{i=1}^n \underline{x}_i / n$$

to the Sebestyen classifier designed with reference set consisting of the \underline{a}_{jk} . The misclassification probability of such a procedure can be expressed as a function of the number of sample patterns as follows.

Let R_j be the region of feature-space associated with class ω_j by the Sebestyen classifier. In the two-class case, regions R_1 and R_2 are defined by Eq. (4-2). Since the distribution of the maximum-likelihood estimate of \underline{a}_{jk} is $n[\underline{a}_{jk}, (\sigma^2/n)I]$ it follows that the misclassification probability associated with subclass ω_{jk} is

$$p_{jk}(n) = 1 - \frac{n^{d/2}}{(2\pi)^{d/2} \sigma^d} \int_{R_j} \exp(-n \|\underline{x} - \underline{a}_{jk}\|^2 / 2\sigma^2) d\underline{x}.$$

The average misclassification probability of the procedure with the estimate defined by Eq. (6-1) as input is given by

$$P_e(n) = (1/cs) \sum_{j=1}^c \sum_{k=1}^s p_{jk}(n)$$

VI.B. A Suboptimal Procedure

The procedure based on submitting an estimate to an available classifier is restrictive in that it assumes constant aircraft orientation with each observation. In

other words it assumes that all noisy patterns belong to the same subclass. In this section we will describe a procedure that utilizes an available classifier which does not make the constant subclass assumption. However, it is an ad hoc procedure that can be expected to perform nearly the same irrespective of the orientation.

Let C be the available classifier, and let $\underline{x}_1, \dots, \underline{x}_N$ be N sequentially arriving patterns. Let c be the number of possible pattern classes. Let $f(n) \leq n \leq N$ where $f(n)$ and n are integers. The suboptimal algorithm, to be called Algorithm M, is as follows.

Submit each of the sequentially arriving patterns to the classifier C . By means of c counters keep track of the number of patterns classified into each class. As soon as $f(n)$ of the patterns are assigned to the same class ω_i , terminate the procedure. The decision of the procedure is taken to be ω_i .

This simple procedure has a number of interesting characteristics. The number of patterns required to arrive at a decision, m , has to lie in the range defined by

$$(6-2) \quad f(n) \leq m \leq [f(n)-1] \cdot c + 1$$

The lower bound is reached when the first $f(n)$ patterns are classified by C into the same class. The upper bound is reached when each of the c classes are represented by $[f(n)-1]$ patterns, and an additional pattern is required to break the tie.

If we choose $f(n) = [(n-1)/c] + 1$ then m lies between $[(n-1)/c] + 1$ and n , as can be verified from inequality (6-2). In order to make $f(n)$ integer-valued we choose $n = c \cdot k + 1$, where $k = 1, 2, 3, \dots$. Note that if $c=2$, the above algorithm chooses that class which was assigned by C to a simple majority of the n patterns.

If we arbitrarily choose $f(n) = (n+1)/2$, which is the value of $[(n-1)/c] + 1$ for $c=2$, then in the general c -class case m lies in the range

$$(n+1)/2 \leq m \leq [(n-1)/2] \cdot c + 1.$$

Note that in this case the upper bound of m is greater than or equal to n .

The actual number of observations, N , required to implement the above algorithm is as follows. If $c=2$, and one desires the procedure to be equivalent to a majority-vote scheme over n decisions, i.e., $f(n) = (n+1)/2$, then N

has to be at least as large as n . In the general c -class case, if one chooses $f(n) = (n+1)/2$, then N has to be at least as large as $(n-1)/2 \cdot c + 1$.

VI.B.1 Misclassification Probability of the Majority-Vote Algorithm

If the class-conditional error probabilities of C are known, the corresponding error probabilities of the above procedure can be expressed as a function of the number of sample patterns n . In particular, we will consider the two-class problem here. It was observed earlier that a choice of $f(n)=(n+1)/2$, n odd, makes the above procedure equivalent to a procedure that takes a majority-vote over n individual decisions. We will examine here the behavior of the average probability of error of the majority-vote algorithm as a function of n .

Consider the two-class problem with classes ω_1 and ω_2 . Let classifier C have the class-conditional probability of error p_i associated with class ω_i . Denote the classifier that is based on $f(n)=(n+1)/2$ as $C(n)$. Let $C(n)$ have class-conditional error probability $p_i(n)$ associated with class ω_i . If $p_i(j;n)$ represents the probability of exactly j out of n decisions of C not being ω_i when ω_i is the true class, then

the probability $p_i(n)$ can be expressed as a sum of its component probabilities as follows:

$$p_i(n) = p_i[f(n);n] + p_i[f(n)+1;n] + \dots + p_i[n;n] .$$

The probability that exactly j out of n decisions of C are not ω_i when ω_i is true is given by the binomial law.

$$\text{That is, } p_i[j;n] = \binom{n}{j} [p_i]^j [1-p_i]^{n-j} .$$

Therefore,

$$p_i(n) = \sum_{j=f(n)}^n \binom{n}{j} [p_i]^j [1-p_i]^{n-j} .$$

Fig. 6-1 illustrates the behavior of $p_i(n)$ as a function of p_i for $c=2$ and varying values of n . This behavior can also be summarized as follows:

If $p_i \leq (1/2)$ then $p_i(n+1) \leq p_i(n)$;

if $p_i = (1/2)$ then $p_i(n+1) = p_i(n) = (1/2)$;

if $p_i \geq (1/2)$ then $p_i(n+1) \geq p_i(n)$.

The average misclassification probability of the majority-vote procedure is given by

$$P_e(n) = \sum_{i=1}^2 P(\omega_i) \cdot p_i(n)$$

or

$$(6-3) \quad P_e(n) = \sum_{i=1}^2 P(\omega_i) \cdot \sum_{j=f(n)}^n \binom{n}{j} (p_i)^j (1-p_i)^{n-j} .$$

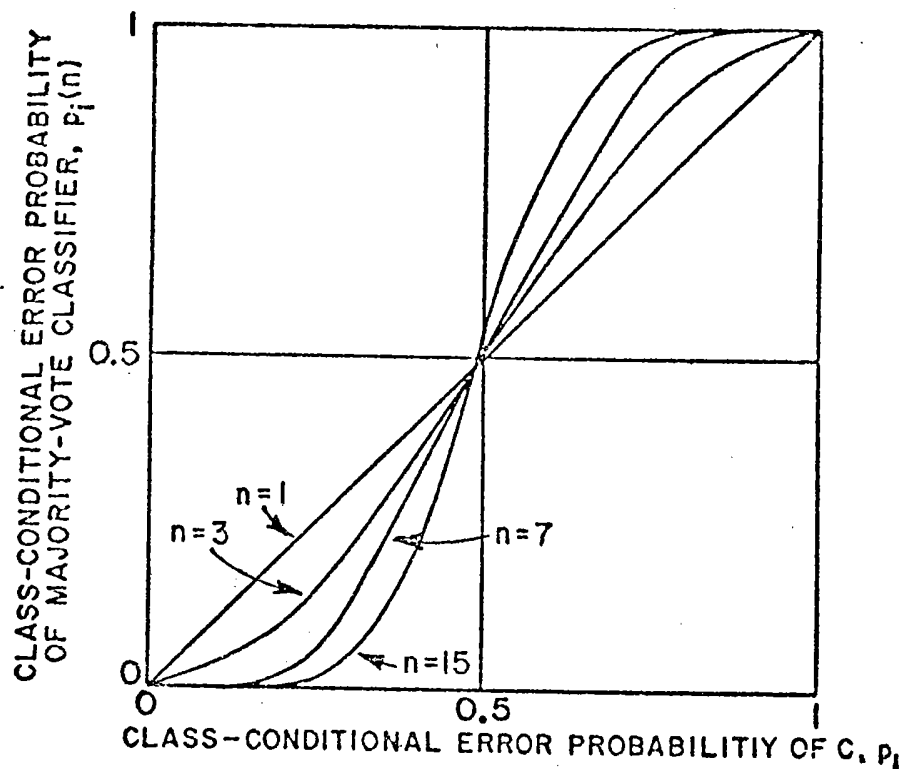


Fig. 6-1. Variation of $p_i(n)$ with p_i . Curves are parametrized on the number of observations.

In any application of a multiple observation procedure one is interested in the number of observations necessary to achieve a specified average misclassification probability. We will show here that this number is dependent on the classification bias of the classifier C , as long as the average misclassification probability of C has some fixed value.

Theorem 6.1: Let C be a classifier for a problem with two equally likely classes such that its average misclassification probability is a , and classification bias is b . For a

constant value of a , the average misclassification probability of a procedure that takes a majority-vote over n decisions of C is a monotonic increasing function of b .

Proof: Let p_1 and p_2 be the class-conditional error probabilities of C . Without loss of generality, assume $p_1 \geq p_2$. Thus, $a = \frac{1}{2}(p_1 + p_2)$ and $b = \frac{1}{2}(p_1 - p_2)$.

Since $p_1 = a + b$, and $p_2 = a - b$, the average misclassification probability of the majority-vote procedure can be written, from Eq. (6-3), as

$$P_e(n) = (1/2) \sum_{j=\frac{n+1}{2}}^n \binom{n}{j} [(a+b)^j (1-a-b)^{n-j} + (a-b)^j (1-a+b)^{n-j}].$$

If the term $[(a+b)^j (1-a-b)^{n-j} + (a-b)^j (1-a+b)^{n-j}]$ is monotonic increasing with b for all possible values of j , then the theorem follows. However, the term is actually monotonic decreasing with b for certain values of a , n , and j , e.g., $a=0.4$, $n=3$, $j=2$. Thus the monotonicity of P_e with b is determined by the total sum.

$$\text{If } B(k; n, p) = \sum_{j=0}^k \binom{n}{j} p^j (1-p)^{n-j},$$

then

$$1 - B(k; n, p) = n \binom{n-1}{k} \int_0^p (1-t)^{n-k-1} dt$$

where the integral is the incomplete beta function (see Feller, 1968).

$$\begin{aligned} n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \\ n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \\ n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \end{aligned}$$

$$\begin{aligned} n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \\ n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \end{aligned}$$

$$k = \frac{n-1}{2}$$

$$\begin{aligned} n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \\ n - \frac{(n-1) \cdot 1}{2} &= \frac{n-1}{2} \end{aligned}$$

Since

$$\begin{aligned} \sum_{j=\frac{n+1}{2}}^n \binom{n}{j} \cdot p^j (1-p)^{n-j} &= 1 - B\left[\frac{n-1}{2}; n, p\right] \\ &= n \binom{n-1}{\frac{n-1}{2}} \int_0^p t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt, \end{aligned}$$

we can write

$$P_e(n) = (1/2) n \binom{n-1}{\frac{n-1}{2}} \left[\int_0^{a+b} t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt + \int_0^{a-b} t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt \right]$$

or,

$$(6-4) \quad P_e(n) = (1/2) n \binom{n-1}{\frac{n-1}{2}} \left[2 \int_0^a t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt + D_1(b) \right]$$

where

$$D_1(b) = \int_a^{a+b} t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt - \int_{a-b}^a t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-3}{2}} dt.$$

$P_e(n)$ is monotone increasing with b if and only if $D_1(b)$ is monotone increasing with b , since the other terms on the right hand side of Eq. (6-4) are independent of b .

Let

$$D_2(b) = \int_a^{a+b} t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-1}{2}} dt - \int_{a-b}^a t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-1}{2}} dt.$$

We will show that $D_2(b)$ is monotone increasing with b . Let

$$f(t) = t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-1}{2}}, \quad t \in (0,1).$$

The first derivative of $f(t)$ with respect to t is positive in the interval $(0, \frac{1}{2})$, and negative in $(\frac{1}{2}, 1)$. That is, $f(t)$ is monotone increasing in $(0, \frac{1}{2})$ and monotone decreasing in $(\frac{1}{2}, 1)$. The function is symmetric about $t = \frac{1}{2}$, and can be represented as shown in Fig. 6-2.

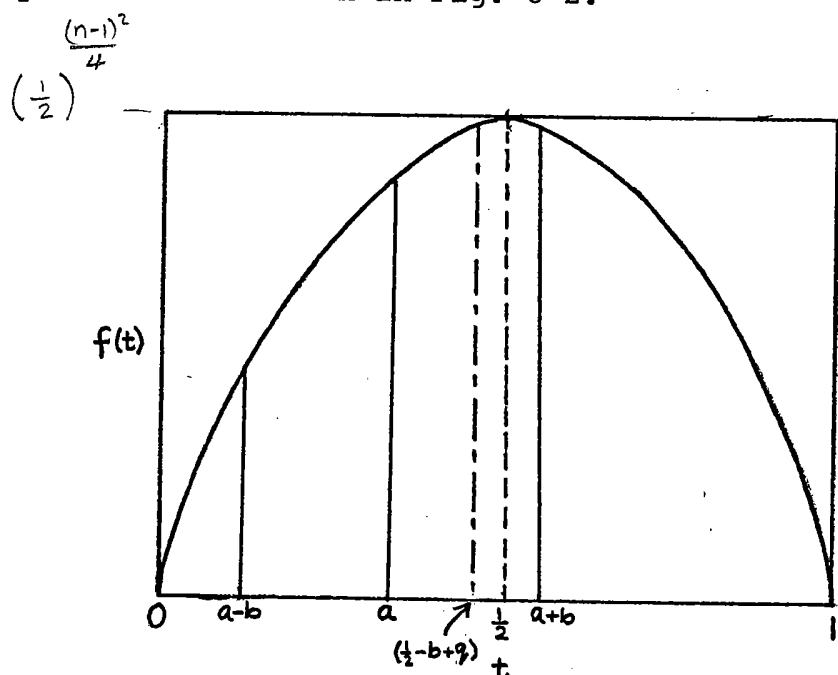


Fig. 6-2. Functional behavior of $f(t) = t^{\frac{n-1}{2}} \cdot (1-t)^{\frac{n-1}{2}}$.

If $a+b \leq \frac{1}{2}$ then $D_2(b)$ is monotone increasing with b , since the first term of $D_2(b)$ increases at a faster rate with b than the second term. If $a+b > \frac{1}{2}$, consider the following. Let $a+q = \frac{1}{2}$. If we show that the area under the curve in interval $(\frac{1}{2}, a+b)$ is larger than the area in interval

$(a-b, a-q)$ it follows that $D_2(b)$ increases monotonically with b . Since the interval $(\frac{1}{2}, a+b)$ is of length $(b-q)$, and $f(t)$ is symmetric about $\frac{1}{2}$, it follows that the area under the curve in interval $(\frac{1}{2}, a+b)$ is equal to that in the interval $(\frac{1}{2}-b+q, \frac{1}{2})$. Due to monotonicity of $f(t)$ in $(0, \frac{1}{2})$, the area in interval $(\frac{1}{2}-b+q, \frac{1}{2}) = (\frac{1}{2}, a+b)$ is greater than the area in the interval $(a-b, a-q)$, and thus $D_2(b)$ increases monotonically with b .

The amount of increase of $D_1(b)$ with b is higher than that of $D_2(b)$, since $D_1(b)$ corresponds to the function $f(t)/(1-t)$ that is a faster increasing function of t than $f(t)$. That is, $D_1(b)$ is monotone increasing with b , or $P_e(n)$ increases monotonically with b . Q.E.D.

Fig. 6-3 illustrates the behavior of $P_e(n)$ with n for the particular situation where $P(\omega_1) = P(\omega_2) = \frac{1}{2}$ and the average misclassification probability of C , $P_e = \frac{1}{2}(p_1 + p_2)$, is 0.3. The curves are parametrized on the classification bias of C .

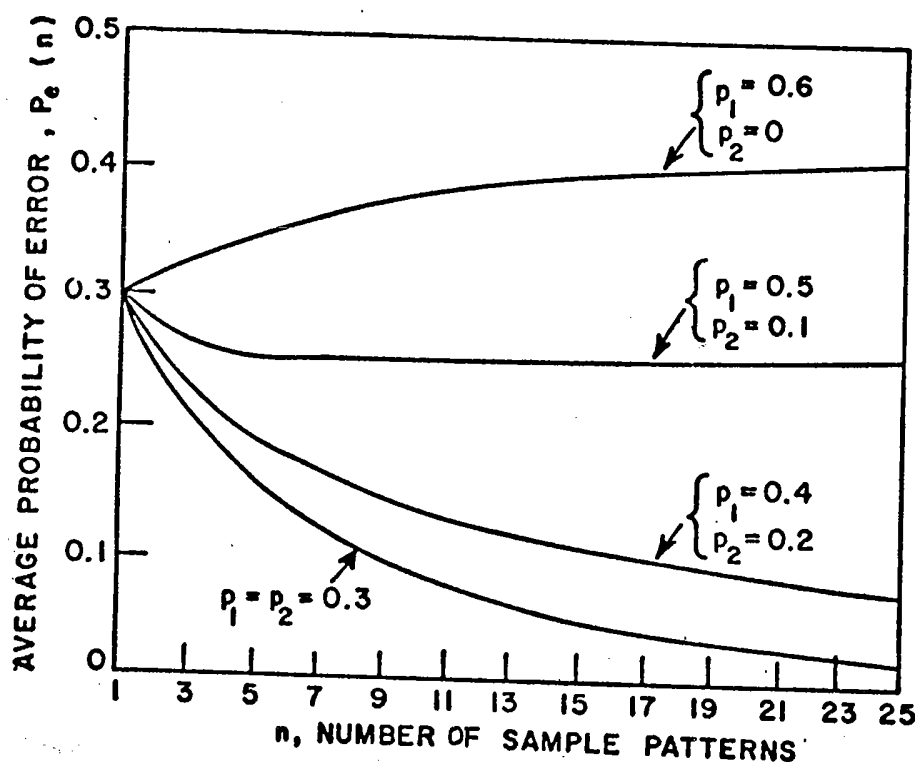


Fig. 6-3. $P_e(n)$ as a function of n , for different values of classifier bias and constant $P_e=0.3$.

Observe that when $p_1 - p_2 < 0.5$, $P_e(n)$ decreases with increasing n . However, when $p_1 - p_2$ exceeds 0.5, $P_e(n)$ increases with n . That is, when the bias of C exceeds 0.25, the performance of the majority-vote procedure deteriorates as more samples are obtained. The increase in $P_e(n)$ with n can be explained as follows. Since $p_1 - p_2 > 0.5$, and $(1/2)(p_1 + p_2) \leq 0.5$, i.e., the average probability of error of any classifier is less than 0.5, it follows that

$$0.5 + p_2 < p_1 \leq 1 - p_2 .$$

In order to satisfy the above inequality, p_1 and p_2 have to lie in the ranges $0.5 \leq p_1 \leq 1$ and $0 \leq p_2 \leq 0.25$. This means that $p_1(n)$ increases with n while $p_2(n)$ decreases with n .

Now if p_1 and p_2 have values such that $p_1(n)$ increases at a faster rate than $p_2(n)$ decreases, then it follows that $P_e(n)$ increases with n . Since $p_1=0.6$ and $p_2=0$ are such values, as can be verified from Fig. 6-1, it follows that $P_e(n)$ increases with n .

VI.C. Performance of Procedures with Radar Aircraft Data

The average misclassification probability of the two fixed sample-size multiple observation procedures described in the previous sections are compared here. The comparison is in terms of their performance with radar aircraft data, when either the nearest-neighbor classifier or a Bayes classifier is utilized. The estimate-based procedure considered is one that submits the sample mean to a given classifier.

Fig. 6-4a shows $P_e(n)$ as a function of n , for each of the four different possibilities. The set of curves

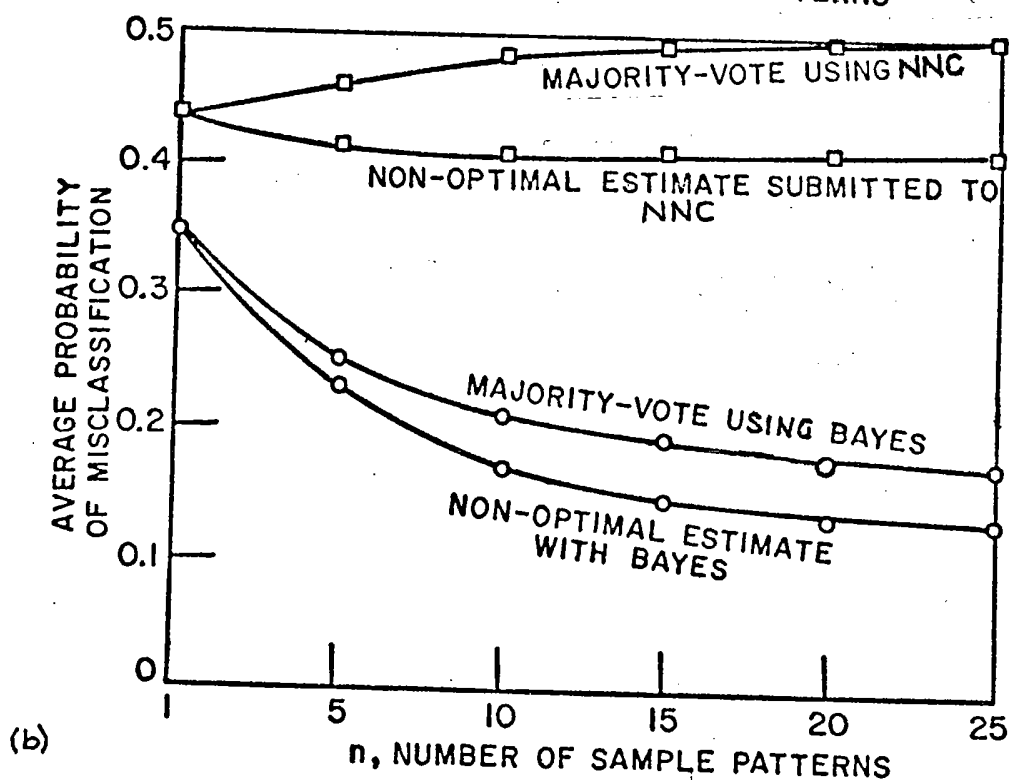
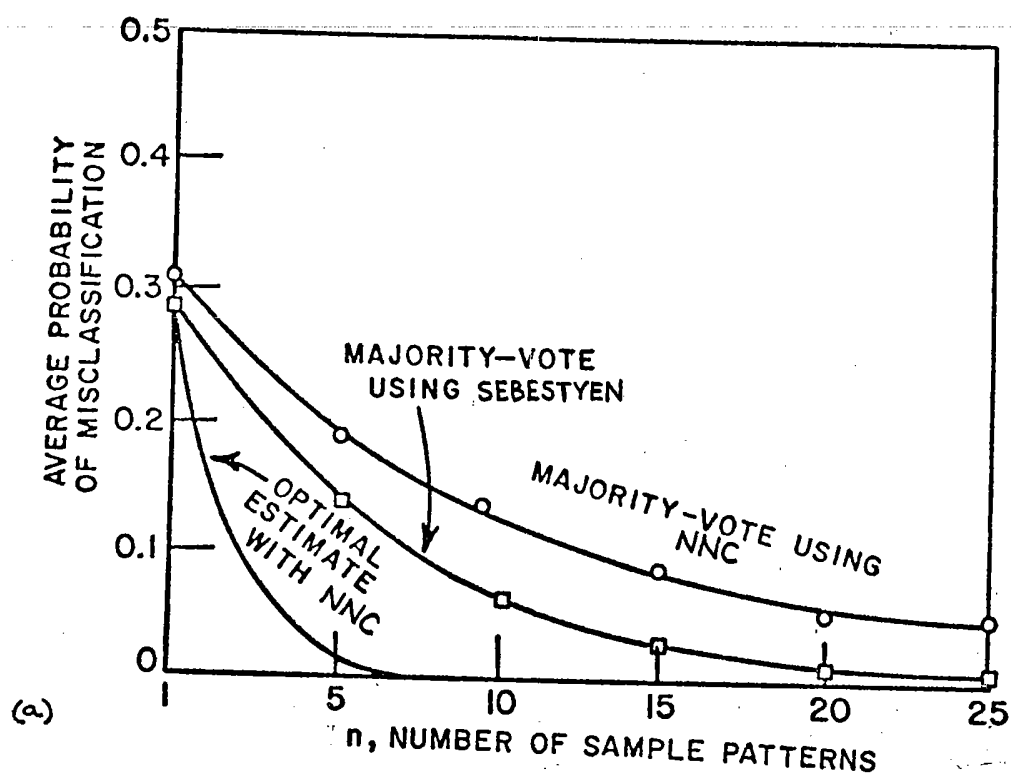


Fig. 6-4. Comparison of performance of fixed sample-size multiple observation procedures with radar aircraft data.

correspond to $p(\underline{x}/\omega_{ij}) \sim n(\underline{a}_{ij}, \sigma^2 \mathbf{I})$, $s=9$, $d=8$, and $\sigma = 1.4$ times the average amplitude return from the F4 and the F-104. The aircraft orientation was assumed to be constant at $(\phi_0, \theta_0) = (0^\circ, 110^\circ)$ as each of the n observations was made. The single observation performance of the Sebestyen and nearest-neighbor classifiers was as follows

$$\text{NNC: } P_e = 0.31, B = 0.15$$

$$\text{SC: } P_e = 0.28, B = 0.015$$

In spite of the fact that the average misclassification probabilities of the Sebestyen and nearest-neighbor classifiers for $n=1$ are near each other, the nearest-neighbor classifier needs considerably more samples than the Sebestyen classifier to reach an error probability of, say, 0.01. This is explained by the high bias of the nearest-neighbor classifier as compared to that of the Sebestyen classifier. The estimate-based procedure is seen to perform considerably better than the majority-vote procedure, i.e., shows rapid decrease of $P_e(n)$ with n . This is expected of the estimate-based procedure since for any given n , the procedure which minimizes misclassification probability is one that submits the maximum-likelihood estimate of the ideal pattern to the Bayes classifier.

Fig. 6-4b corresponds to the situation where class-conditional densities are Rician mixtures, with $d=4$. In this case the single-sample error probabilities of the classifiers were as follows:

NNC: $P_e=0.44$, $B=0.32$

Bayes (based on Rician): $P_e=0.36$, $B=0.08$.

In this case, the estimate-based procedure does not show marked improvement over the majority-vote procedure, a fact explained by the suboptimal nature of the estimate.

VI.D. Sequential Classification

If class-conditional probability densities are known, then standard sequential hypothesis testing procedures can be used in the design of sequential procedures for the multiple observation problem. In the two-class case, assuming independence of observations, the Wald's sequential probability ratio test (SPRT) can be used as follows. Let the class-conditional densities be mixtures of s equally likely Gaussian densities with a common variance σ^2 . The likelihood ratio at the n^{th} stage of the sequential process is defined by

$$L_n = \frac{p(\underline{x}_1, \dots, \underline{x}_n / \omega_1)}{p(\underline{x}_1, \dots, \underline{x}_n / \omega_2)}$$

$$= \frac{\prod_{i=1}^n \sum_{k=1}^S \exp \left[- (1/2\sigma^2) (\underline{x}_i - \underline{a}_{1k})^t (\underline{x}_i - \underline{a}_{1k}) \right]}{\prod_{i=1}^n \sum_{k=1}^S \exp \left[- (1/2\sigma^2) (\underline{x}_i - \underline{a}_{2k})^t (\underline{x}_i - \underline{a}_{2k}) \right]} .$$

If the desired type 1 and type 2 error probabilities are α and β , then the procedure is to

choose class ω_1 if $L_n > \frac{1-\beta}{\alpha}$,

choose class ω_2 if $L_n < \frac{\beta}{1-\alpha}$,

and obtain the next sample if neither condition is true.

The thresholds $\frac{1-\beta}{\alpha}$ and $\frac{\beta}{1-\alpha}$ are approximations to the optimal stopping boundaries that minimize the expected number of observations. In practice, the approximate thresholds are known to perform well. Further discussion of stopping boundaries for the sequential probability ratio test can be found in Wald (1947).

VI.D.1. Sequential Classification Using the Bayes Classifier

The single observation Bayes classifier utilizes the class a priori probabilities and class-conditional densities in order to compute the posteriori probabilities of classes when an observed pattern is given as input. Consider the following procedure: when the i^{th} observation

is made, utilize the class a posteriori probabilities from the $(i-1)^{\text{th}}$ observation as the class a priori probabilities in order to compute the i^{th} a posteriori probabilities; terminate the procedure if one of the a posteriori probabilities exceeds a specified threshold; if not, make the $(i+1)^{\text{th}}$ observation. This procedure has been referred to as the Bayes sequential procedure in the literature (Roberts and Mullis, 1970). Fig. 6-5 is a schematic illustration of the Bayes sequential procedure applied to the radar aircraft identification problem.

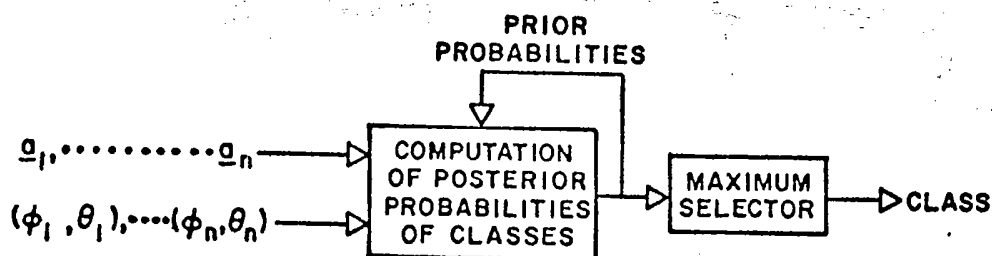


Fig. 6-5. The Bayes Sequential Procedure

We will show how the thresholds for such a procedure can be directly derived from the Wald's sequential probability ratio test in the two-class case, and from the multiple hypothesis testing procedure due to Armitage in the c -class case. We will then evaluate the performance of

the two-class procedure for the radar aircraft identification problem when the Bayes classifier has the form of the Sebestyen classifier, i.e., the subclass distribution is Gaussian.

Consider first the two-class case. The Wald's sequential probability ratio test chooses class ω_1 if

$$\frac{p(\underline{x}_1, \dots, \underline{x}_n / \omega_1)}{p(\underline{x}_1, \dots, \underline{x}_n / \omega_2)} > \frac{1-\beta}{\alpha}.$$

This is equivalent to choosing class ω_1 if

$$(6-5) \quad \frac{p(\underline{x}_1, \dots, \underline{x}_n / \omega_1) \cdot P(\omega_1)}{\sum_{i=1}^2 p(\underline{x}_1, \dots, \underline{x}_n / \omega_i) \cdot P(\omega_i)} > \frac{(1-\beta) \cdot P(\omega_1)}{\alpha \cdot P(\omega_2) + (1-\beta) \cdot P(\omega_1)}$$

Inequality (6-5) can be written as

$$(6-6) \quad P(\omega_1 / \underline{x}_1, \dots, \underline{x}_n) > \frac{(1-\beta) \cdot P(\omega_1)}{\alpha \cdot P(\omega_2) + (1-\beta) \cdot P(\omega_1)},$$

where $P(\omega_1 / \underline{x}_1, \dots, \underline{x}_n)$ is the a posteriori probability of class ω_1 after n observations have been made. The Bayes sequential procedure computes the posterior probability, $P(\omega_1 / \underline{x}_1, \dots, \underline{x}_n)$, by utilizing $P(\omega_1 / \underline{x}_1, \dots, \underline{x}_{n-1})$ as the prior probability.

We have thus obtained the thresholds in terms of the class a priori probabilities, and desired class-conditional error probabilities. It should be noted that the thresholds are only as good as the initial approximations to the optimal stopping boundaries that we started out with. Note that the test we have obtained is algebraically identical to the Wald's sequential probabilities ratio test, i.e., the two tests always choose the same class irrespective of the class a priori probabilities. The only advantage to expressing the Wald's test as we have done is that the Bayes classifier designed for the single observation problem can now be implemented as a sequential procedure that is identical to the Wald's test. This eliminates the need for two separate procedures for the single and multiple observation problems.

In the c-class case, there is no single multiple sequential hypothesis testing procedure that requires the fewest number of observations to reach the desired misclassification levels for arbitrary class-conditional

densities. Among various suboptimal multiple sequential hypothesis testing procedures are the generalized sequential probability ratio test (GSPRT) and the test due to Armitage (Fu, 1968). We will describe here the Armitage test, and obtain the thresholds for the Bayes sequential procedure from those of the Armitage test.

The Armitage procedure involves performing likelihood ratio tests between every pair of hypotheses and accepting a hypothesis if it is accepted over all others. If no such hypothesis exists, the next observation is made. The formulation of the test is as follows.

Accept hypothesis H_i if all of the following $(c-1)$ inequalities are true, for $i=1, \dots, (c-1)$.

$$(6-8) \quad \frac{p(\underline{x}_1, \dots, \underline{x}_n / \omega_i)}{p(\underline{x}_1, \dots, \underline{x}_n / \omega_j)} > A_{ij} \quad j=1, \dots, c, j \neq i ;$$

Accept hypothesis H_c if all of the following inequalities are true,

$$\frac{p(\underline{x}_1, \dots, \underline{x}_n / \omega_i)}{p(\underline{x}_1, \dots, \underline{x}_n / \omega_c)} < (1/A_{ci}) ;$$

where

$$A_{ij} = \left[1 - \sum_{\substack{k=1 \\ k \neq j}}^c p_{kj} \right] / p_{ij}$$

and p_{ij} is the desired probability of accepting H_i when H_j is true, i.e., the transitional error probability required of the procedure. In the case where $c=2$, the Armitage procedure can be verified to be identical to Wald's sequential probability ratio test.

The sequential procedure defined above has the property that the probability of no decision being reached at the n th stage tends to zero as n increases indefinitely (Armitage, 1950). By choosing the A_{ij} sufficiently large we can make the probability of correct recognition, when any of the H_i is true, as large as we wish. However, the test is conservative in the sense that the true probabilities tend to be lower than the values incorporated into the stopping boundaries.

By multiplying each of the inequalities of (6-8) by $1/P(\omega_j)$, changing inequality signs by inverting, and by adding each of the $(c-1)$ inequalities, the Armitage procedure is equivalent to accepting hypothesis H_i if

$$\frac{\sum_{\substack{j=1 \\ j \neq i}}^c p(\underline{x}_1, \dots, \underline{x}_n / \omega_j) \cdot P(\omega_j)}{p(\underline{x}_1, \dots, \underline{x}_n / \omega_i)} < \sum_{\substack{j=1 \\ j \neq i}}^c \frac{P(\omega_j)}{A_{ij}} .$$

That is,

$$(6-9) \quad P(\omega_i / \underline{x}_i, \dots, \underline{x}_n) > \frac{P(\omega_i)}{P(\omega_i) + \sum_{\substack{j=1 \\ j \neq i}}^c [P(\omega_j) / A_{ij}]}$$

The right hand side of inequality (6-9) is the desired threshold. A case of interest is when class a priori probabilities are equal, and the desired transitional probabilities are all equal to α . If these values are substituted into (6-9), the threshold reduces to $1-(c-1)\alpha$. That is, the Bayes sequential procedure in this case is to continue observations until one of the elements of the Bayesian probability vector $P_n = [P_n(\omega_1), \dots, P_n(\omega_c)]$ exceeds $1-(c-1)\alpha$.

The Bayes sequential procedure in the c -class case is identical to the Armitage procedure for sequential multi-hypothesis testing. In other words, the two procedures choose the same class at the same stage of the sequential process. The identity of the two procedures can be verified to be true for arbitrary a priori probabilities of class membership.

As in the two-class case, the advantage of expressing the sequential procedure as we have done lies in the fact that it enables one to use the Bayes classifier in a

sequential mode. Since the Armitage procedure is not guaranteed to require the minimum number of observations for achieving a specified misclassification probability, no claim as to optimality of the procedure can be made. In the following section we will experimentally evaluate the performance of the Bayes sequential procedure with radar aircraft data. The situation to be considered is one where the Bayes classifier has the form of the Sebestyen classifier.

VI.D.2. Application of Sequential Procedure to Radar Aircraft Identification

The a posteriori probability of class ω_i as computed by the Sebestyen classifier at the n^{th} stage of the sequential process is

$$P_n(\omega_i) = \frac{\sum_{(\emptyset, \theta) \in R(\emptyset_n, \theta_n)} \exp \left[-\| \underline{x}_n - \underline{f}_i(\emptyset, \theta) \|^2 / 2\sigma^2 \right] \cdot P_{n-1}(\omega_i)}{\sum_{j=1}^C \sum_{(\theta, \emptyset) \in R(\emptyset_n, \theta_n)} \exp \left[-\| \underline{x}_n - \underline{f}_j(\theta, \emptyset) \|^2 / 2\sigma^2 \right] \cdot P_{n-1}(\omega_j)}$$

where \underline{x}_n is the observed pattern at the n^{th} stage, (\emptyset_n, θ_n) is the corresponding observed orientation, $R(\emptyset_n, \theta_n)$ is the set of orientations that have a non-zero probability of being

the true observation, $f_1(\emptyset, \theta)$ is the reference pattern corresponding to (\emptyset, θ) and σ^2 is the noise variance.

The application of the above procedure to the radar aircraft identification problem yielded the results shown in Tables 3 and 4. In obtaining these results, amplitude and phase at four frequencies were used to construct the pattern vector. As a consequence the subclass distributions were multivariate Gaussian with $d=8$. Hypothesis H_0 corresponds to the F-104, and hypothesis H_1 to the F-4. The class a priori probabilities were assumed to be equal, and $R(\emptyset_n, \theta_n)$ to be uniform over nine discrete orientations centered at $(0^\circ, 30^\circ)$. The thresholds used to terminate the sequential procedure were the right hand sides of inequalities (6-6) and (6-7). The values of class-conditional error probabilities used to compute the thresholds are given in the first column of each table, i.e., α and β .

The experimental error probabilities tended to be lower than the error probabilities used to compute the thresholds. This agrees with what is known about the Wald's sequential probability ratio test, i.e., the thresholds $\frac{1-\beta}{\alpha}$ and $\frac{\alpha}{1-\beta}$ yield error probabilities α' and β' such that

TABLE 3

Noise Standard Deviation = .5

Error probabilities used to compute thresholds	Experimental error probability		Average number of observations per decision	
	H_0	H_1	H_0	H_1
$\alpha = \beta$				
.1	.009	.018	1.22	1.30
.05	.004	.018	1.35	1.42
.01	.000	.004	1.67	1.73

TABLE 4

Noise Standard Deviation = 1.0

Error probabilities used to compute thresholds	Experimental error probability		Average number of observations per per decision	
	H_0	H_1	H_0	H_1
$\alpha = \beta$				
.1	.044	.094	2.82	3.00
.05	.027	.071	3.76	3.81
.01	.004	.051	5.74	6.90

Note: Hypotheses H_0 and H_1 correspond to aircraft types F-104 and F-4 respectively.

$(\alpha' + \beta') \leq (\alpha + \beta)$. An exception was the case where $\sigma = 1$. and $\alpha = \beta = .01$. This is a situation where the noise standard deviation is larger than the average interclass separation of ideal patterns, and the desired error probability is very small. Considering this to be an unrealistic case, it can be said that the Bayes sequential procedure with thresholds defined by Eq. (6-9) is a conservative test with radar aircraft data. The tables show that when $\sigma = 0.5$, or the noise standard deviation is equal to half the average amplitude return, the average number of observations required to reach a misclassification probability of .01 is less than two. As expected this number increases as the noise standard deviation is increased. We conclude this section by observing that when data from the F-4 and the F-104 is used, with realistic values of σ ($\sigma < 0.5$), arbitrarily small misclassification levels can be reached with few observations.

CHAPTER VII

SUMMARY AND CONCLUSIONS

VII.A. Summary

Aspects of the problem of classifier design for a type of pattern recognition problem have been investigated. The classes are characterized by highly convoluted, "ideal" data surfaces in feature-space. The measured pattern is considered to be a "noisy" version of an "ideal" pattern that lies on one of the data surfaces. The problem of radar aircraft identification has been considered to be an example of such a problem, where the data surfaces are functions of aircraft orientation with respect to the radar.

A formulation of the classification problem has been made in terms of Bayes decision theory. This includes statistical models of orientation uncertainty, and uncertainty associated with feature measurements for any given

orientation value. Based on this formulation the Bayes classifier for the problem has been derived.

Considering the particular case where orientation uncertainty is modeled by a uniform distribution, and feature measurement uncertainty, or "noise", is modeled by a normal distribution, some problems concerning the implementation and evaluation of the Bayes classifier have been investigated. These problems have included:

- (i) efficient computer implementation of the Bayes classifier;
- (ii) evaluation of misclassification probability of the Bayes classifier;
- (iii) implementation of the Bayes classifier when a parameter of feature measurement uncertainty, viz., standard deviation of the normal distribution, is inexactly determined;
- (iv) effect of increasing dimensionality on misclassification probability at individual orientation possibilities;
- (v) utilization of the Bayes classifier in a sequential multiple observation procedure.

The Bayes classifier for the problem where class-conditional probability distributions are Gaussian mixtures

was first implemented in a practical pattern recognition problem by Sebestyen (1962). Since the nearest-neighbor classifier with a reference set corresponding to discrete orientation possibilities has been used with success in the radar aircraft identification problem, properties of the Sebestyen and nearest-neighbor classifiers have been compared here.

A relationship between the Sebestyen and nearest-neighbor classifiers has been obtained by showing that the former is obtained when discriminant functions are based on mixtures of distributions, with the mixing being over orientation possibilities, and the latter is obtained when multiple discriminant functions corresponding to each class and orientation possibility are used to determine the subclass and corresponding class. Necessary and sufficient conditions on reference pattern configuration for the identity of decision surfaces of the two classifiers have been obtained. Based on these conditions, an algorithm for determining the identity of decision surfaces has been provided.

The Sebestyen and nearest-neighbor classifiers have been compared in terms of standard classifier evaluation

criteria, viz., misclassification probability and computational economy, and in terms of criteria which were thought to be appropriate for radar aircraft identification, viz., extensibility of the classifier to the multiple observation problem, and classifier robustness. Misclassification probability, classification bias, and robustness were experimentally compared with multivariate data from the radar aircraft identification problem. Comparison of computational economy was made in terms of number of computation steps, and computation time on a typical general purpose computer. The effect of increase in dimensionality on both classifiers were examined analytically in restricted situations. It was observed that the misclassification probability of both classifiers for particular orientations can increase with dimensionality.

VII.B. Specific Contributions and Conclusions

The discriminant functions of the Sebestyen classifier and the Bayes classifier based on Rician mixtures are sums of exponentials. A computational procedure for comparing such discriminant functions has been obtained which overcomes the problem of overflow/underflow in a

computer implementation. The procedure tends to need to compute far fewer exponentials than the direct procedure; in the best case every exponentiation is replaced by a single exponentiation and two comparisons. Thus the procedure is, in general, computationally faster than the direct procedure. If the procedure is used to implement the Sebestyen classifier, the probability of needing to compute any exponentials by the procedure is small when the standard deviation of feature measurement uncertainty is low compared to the inter-class separation of patterns.

Bounds on the misclassification probability of the Sebestyen classifier, when it is the Bayes classifier, were obtained in terms of the inter-class separation of patterns. The utility of the Bhattacharyya coefficient in approximating the misclassification probability of the Sebestyen classifier has been shown. An analysis of the dependency of the Bhattacharyya coefficient on the value of standard deviation of feature measurement uncertainty has been made. The utility of this analysis is that the behavior of the Sebestyen classifier with increasing feature measurement uncertainty can be predicted without having to consider many possible values of the standard deviation. For

example, in the radar aircraft identification problem the performance of a classifier at different noise levels affecting signals is of interest. In such a situation one can specify the significant noise levels to be considered by an inspection of the configuration of reference patterns.

For an arbitrary set of reference patterns, the Sebestyen and nearest-neighbor classifiers have been shown to be related. A test to determine whether the two classifiers result in identical decision surfaces has been provided. When performed prior to classifier implementation, such a test can lead to more efficient classifier operation.

The variation of the decision boundary of the Sebestyen classifier with an assumed value of noise standard deviation has been analytically examined in the linearly separable univariate case. It was found that an underestimated value of noise standard deviation leads to a smaller change from the optimal decision boundary than an overestimated value. In the multivariate case, the effect of misclassification probability of the Sebestyen classifier when the assumption of noise standard deviation

is incorrect was experimentally examined. Here again, it was found that an underestimated value leads to a smaller increase in the misclassification probability than an overestimated value.

When a classifier designed on the basis of inexact orientation knowledge is utilized in a situation where the possible orientations form a proper subset of the orientations considered by the classifier, an increase in feature space dimensionality can lead to an increase in the relevant misclassification probability. Considering the Sebestyen classifier, the following effect of increased feature-space dimensionality has been demonstrated. If the reference pattern configurations in the original and updated feature-spaces satisfy specified geometric constraints, then certain types of features lead to an increase in the relevant misclassification probability. In the case of the nearest-neighbor classifier a similar increase in misclassification probability with dimensionality occurs. However, the causes for the effect need not be the same as those for the Sebestyen classifier. The utility of the analysis of dimensionality is primarily to serve as a warning to the user that certain counter-intuitive phenomenon can occur

even when reference patterns are ideal and the classifier is the Bayes classifier for the problem.

As a consequence of the comparison of properties of the Sebestyen and nearest-neighbor classifiers, the following conclusions can be made. In terms of economy of implementation the nearest-neighbor classifier tends to be more efficient than the Sebestyen classifier. However, the difference in computation time on a typical general-purpose computer is insignificant when dimensionality is higher than a moderate value (about four). This conclusion is independent of the number of classes, and the number of discrete orientation possibilities. In terms of experimentation with data from the radar aircraft identification problem, it was found that the Sebestyen and nearest-neighbor classifier misclassification probabilities tend to be the same, whereas the difference between their class-conditional error probabilities, termed classification bias, tends to be different. The lower classification bias of the Sebestyen classifier was found to be of significance in the multiple observation problem. When class-conditional densities are Rician mixtures, it was found that the performance of the Sebestyen classifier, with a reasonable

assumption of noise variance, is similar to that of the nearest-neighbor classifier. This latter fact helps demonstrate the robustness of the Sebestyen classifier with one type of non-Gaussian noise that is encountered in radar aircraft identification.

In the multiple observation classification problem, using a fixed number of observations, both the Sebestyen and nearest-neighbor classifiers perform similarly when an estimate of the ideal pattern is submitted to the classifier. However, if aircraft orientation changes with each new observation, an estimate-based multiple observation procedure is not feasible; since each of the reference patterns can correspond to a different subclass, and the reference sets of patterns for each observation can be different, even disjoint. A procedure that takes a majority-vote over n individual decisions of a given classifier has been analyzed. The misclassification probability of such a procedure, in the general c -class case, has been expressed as a function of the transition error probability of the given classifier. In the two-class case it has been shown that for a fixed average misclassification probability of a

classifier, the performance of the majority-vote procedure deteriorates with increasing classification bias.

Since the nearest-neighbor classifier tends to have a higher classification bias than the Sebestyen classifier in radar aircraft identification, the majority-vote procedure based on the nearest-neighbor classifier is more error-prone than the procedure based on the Sebestyen classifier. Consequently, the former procedure tends to need more observations to reach a specified misclassification probability.

VII.C. Future Research

In this section we will indicate some of the open problems concerning classifier design and evaluation that follow from this investigation. These problems refer to the situation where classes are characterized by subclasses such that every subclass corresponds to a particular value in a domain set. Further it is assumed that the subclasses are situated in feature-space such that the classification problem is not adequately solved by a linear discriminant technique.

One of the problems that arises in implementing a classifier for highly convoluted data surfaces is that of

determining the number of discrete values on the data surfaces. That is, given that the orientation possibilities lie in a particular range of values, how many discrete orientation possibilities in this range ought to be considered by a classifier? It is obviously related to the rate of variation of the data surfaces over the range of orientation possibilities. Since radar aircraft data show periodic variations with orientation coordinates, multi-dimensional Fourier techniques are a potential tool in determining the frequency components of the data. Based on a knowledge of the maximum frequency components along each feature, the sampling rate for the data can be determined.

A question related to the above is that of removing redundant reference patterns from any given classifier which requires the explicit storage of reference patterns. For a classifier requiring the storage of all reference patterns, it is possible that many of the reference patterns do not contribute appreciably to the determination of the decision surface. Some work on this problem, for the nearest-neighbor classifier, can be found in the literature (Batchelor, 1974). However, practical algorithms for isolating and removing redundant reference patterns need to be developed.

There exists little analytical work in the area of robustness and stability of classifiers in the multivariate case. It is not clear whether the Sebestyen classifier is more robust than the nearest-neighbor classifier in the presence of non-Gaussian noise other than the type considered here. One of the problems in analyzing classifier robustness is that of characterizing non-normal distributions that are relevant to an application and are at the same time analytically manageable. Lachenbruch, et al. (1973), who showed that the linear discriminant function is greatly affected by non-normality, suggest three types of distributions that are generated from the normal distribution by using nonlinear transformations. A study of classifier performance with such distributions could provide better insight into the problem of relative robustness of the Sebestyen and nearest-neighbor classifiers.

The question of classifier stability refers to the magnitude of change in classifier performance with small changes in the values of stored parameters. Preliminary analysis indicates that the Sebestyen decision boundary is less susceptible to large changes than the nearest-neighbor decision boundary when the reference patterns are perturbed.

A more conclusive analysis of the dependence of misclassification probability on stored pattern values needs to be made. The suggestion that one ought to consider changes in misclassification probability, rather than changes in decision boundaries, is based on two arguments. The first is the absence of an adequate quantitative measure for the difference between two arbitrary decision surfaces. The second is the familiar problem of high dimensionality, which makes the analysis of decision surfaces an extremely complex task.

Throughout our work we assumed that ideal reference patterns corresponding to different values of a function domain are available and that a reasonable estimate of standard deviation of feature measurement uncertainty is also available. In some problem situations this may not be the case. All that may be available is a set of labelled sample patterns corresponding to different classes and values in the function domain. In such a case, estimates of the ideal reference patterns and subclass variance can be utilized in implementing the Sebestyen classifier. Theoretical analysis of the misclassification probability of such an estimate-based procedure does not appear to have been made.

APPENDIX A

Behavior of Bhattacharyya Coefficient with Noise Standard Deviation

Procedure for Determining σ_0

In Section III.D. it was shown that if $(\lambda_2/\lambda_1) > 1.86$ then $\frac{d\rho}{d\sigma}$ is unimodal with a maximum value at $\sigma = \sigma_0$. The value of σ_0 can be determined by a procedure to be developed here.

When $f(\sigma^2) = f_1(\sigma^2) + f_2(\sigma^2)$, σ_0^2 is determined by the equation

$$(A-1) \sum_{i=1}^2 A_i \exp(-A_i/8 \sigma_0^2) \{2A_i/\sigma_0^2 - 3\} = 0$$

If ξ_i is the value of σ^2 at which $f_i'(\sigma^2)$ is maximum, i.e., $\xi_i = (2/3)A_i$, Eq. (A-1) reduces to

$$\sigma_0^2 = (5.33/(\xi_1 - \xi_2)) \cdot \ln \left\{ \frac{\xi_1^2 - \xi_1 \cdot \sigma_0^2}{\xi_2 \cdot \sigma_0^2 - \xi_2^2} \right\}.$$

As long as $(A_2/A_1) < 1.86$, the solution to this equation is a unique value in the interval (ξ_1, ξ_2) , which can be determined numerically by an iterative algorithm. Since ρ has s additive components, the value of σ_0^2 at which $d\rho/d\sigma$ is maximum can be determined by $(s-1)$ applications of the algorithm.

Existence of multiple zeros of $\frac{d^2 \rho}{d\sigma^2}$ when $(\lambda_2/\lambda_1) > 1.86$

The case where $(\frac{\lambda_2}{\lambda_1}) > 1.86$ can be considered in terms of two intervals. In the following it is assumed that $f(\sigma^2) = f_1(\sigma^2) + f_2(\sigma^2)$ where $f_i(\sigma^2) = \exp(-A_i/8\sigma^2)$. If $0.66A_1 < 0.27A_2$ and $0.66A_2 < 1.23A_1$, that is $2.44 > (A_2/A_1) > 1.86$, then $f''(\sigma^2)$ has at most one zero in the interval $(0.66A_1, 1.23A_1)$ and an arbitrary number of zeros in $(1.23A_1, 0.66A_2)$. If $1.23A_1 > 0.27A_2 > 0.66A_1$, that is $4.55 > (A_2/A_1) > 2.44$, then $f''(\sigma^2)$ has an arbitrary number of zeros in $(0.66A_1, 0.27A_2)$ and $(1.23A_1, 0.66A_2)$, and at most one zero in $(0.27A_2, 1.23A_1)$. Finally, if $(A_2/A_1) > 4.55$, an arbitrary number of zeros are possible in $(0.66A_1, 0.66A_2)$.

APPENDIX B

In this appendix the proof of a theorem pertaining to the identity of the Sebestyen and nearest-neighbor classifiers is given. The proof of the theorem is based on Lemma 4.2, which will be considered first.

Lemma 4.2: If $\underline{a}_1, \underline{a}'_1, \underline{a}_2, \underline{a}'_2$ are four distinct points in d -dimensional feature-space, then there exist no two points $\underline{\lambda}$ and $\underline{\gamma}$ which simultaneously satisfy the following criteria:

$$(i) \quad \|\underline{\lambda} - \underline{a}_1\|^2 = \|\underline{\lambda} - \underline{a}_2\|^2$$

$$(ii) \quad \|\underline{\lambda} - \underline{a}'_1\|^2 = \|\underline{\lambda} - \underline{a}'_2\|^2$$

$$(iii) \quad \|\underline{\gamma} - \underline{a}_1\|^2 > \|\underline{\gamma} - \underline{a}_2\|^2$$

$$(iv) \quad \|\underline{\gamma} - \underline{a}'_1\|^2 < \|\underline{\gamma} - \underline{a}'_2\|^2$$

Proof: The proof is by contradiction. The four criteria can be rewritten as

$$(i) \quad \|\underline{a}_1\|^2 - \|\underline{a}_2\|^2 = 2\underline{\lambda}^t(\underline{a}_1 - \underline{a}_2)$$

$$(ii) \quad \|\underline{a}'_1\|^2 - \|\underline{a}'_2\|^2 = 2\underline{\lambda}^t(\underline{a}'_1 - \underline{a}'_2)$$

$$(iii) \quad \|\underline{a}_1\|^2 - \|\underline{a}_2\|^2 > 2\underline{\gamma}^t(\underline{a}_1 - \underline{a}_2)$$

$$(iv) \quad \|\underline{a}'_1\|^2 - \|\underline{a}'_2\|^2 < 2\underline{\gamma}^t(\underline{a}'_1 - \underline{a}'_2).$$

Criteria (i) and (iii) imply that

$$(v) \quad \underline{\lambda}^t(\underline{a}_1 - \underline{a}_2) > \underline{\gamma}^t(\underline{a}_1 - \underline{a}_2).$$

Criteria (ii) and (iv) imply that

$$(vi) \quad \underline{\lambda}^t(\underline{a}_1' - \underline{a}_2') < \underline{\gamma}^t(\underline{a}_1' - \underline{a}_2') .$$

There exist no two vectors $\underline{\lambda}$ and $\underline{\gamma}$ such that inequalities (v) and (vi) are simultaneously satisfied.

Q.E.D.

Theorem 4.1: Let $\{\underline{a}_{1j}\} \in \omega_1$ and $\{\underline{a}_{2j}\} \in \omega_2$ be ideal reference patterns. The nearest-neighbor and Sebestyen classifiers are identical if for every $\underline{\lambda}$ on the nearest-neighbor decision surface, and any \underline{a}_{1j} distant k_j from $\underline{\lambda}$, there exists a unique \underline{a}_{2j} distant k_j from $\underline{\lambda}$.

Proof: Let D_1 and D_2 be the sets of points on the nearest-neighbor and Sebestyen decision surfaces. By definition, every point $\underline{\lambda} \in D_1$ satisfies the following

$$(B-1) \quad \min_j \{ \|\underline{\lambda} - \underline{a}_{1j}\|^2 \} = \min_j \{ \|\underline{\lambda} - \underline{a}_{2j}\|^2 \} .$$

According to the condition of the theorem, for any \underline{a}_{1j} there exists an \underline{a}_{2j} such that

$$\|\underline{\lambda} - \underline{a}_{1j}\|^2 = \|\underline{\lambda} - \underline{a}_{2j}\|^2 = k_j .$$

That is,

$$\exp\{-\|\underline{\lambda} - \underline{a}_{1j}\|^2 / 2\sigma^2\} - \exp\{-\|\underline{\lambda} - \underline{a}_{2j}\|^2 / 2\sigma^2\} = 0 .$$

As a consequence, $\underline{\lambda}$ satisfies the equation

$$(B-2) \quad \sum_{j=1}^S \left\{ \exp[-\|\underline{x} - \underline{a}_{1j}\|^2 / 2\sigma^2] - \exp[-\|\underline{x} - \underline{a}_{2j}\|^2 / 2\sigma^2] \right\} = 0 .$$

Thus $\underline{\lambda}$ is a point on the Sebestyen decision surface, or

$\underline{\lambda} \in D_2$. Since this is true for every $\underline{\lambda} \in D_1$, we can write

$$D_1 \subseteq D_2.$$

It needs to be shown that $D_1 \equiv D_2$. We will let $\underline{y} \in D_2 - D_1$, and show that there can be no such \underline{y} . As a consequence it will be concluded that every point in D_2 has to be in D_1 .

Since $\underline{y} \in D_2$, it follows that Eq. (B-2) is satisfied if we replace \underline{x} by \underline{y} . Let

$$\begin{aligned} \min_j \{ \|\underline{y} - \underline{a}_{1j}\|^2 \} &= \|\underline{y} - \underline{a}_1\|^2, \\ \min_j \{ \|\underline{y} - \underline{a}_{2j}\|^2 \} &= \|\underline{y} - \underline{a}_2\|^2, \end{aligned}$$

where $\underline{a}_1 \in \omega_1$ and $\underline{a}_2 \in \omega_2$. Since $\underline{y} \notin D_1$ it follows that

$$\|\underline{y} - \underline{a}_1\|^2 \neq \|\underline{y} - \underline{a}_2\|^2.$$

Without loss of generality let

$$\|\underline{y} - \underline{a}_1\|^2 > \|\underline{y} - \underline{a}_2\|^2.$$

In order that Eq. (B-2) is satisfied with \underline{y} in place of \underline{x} , there must be an $\underline{a}'_1 \in \omega_1$ and $\underline{a}'_2 \in \omega_2$ such that

$$\|\underline{y} - \underline{a}'_1\|^2 < \|\underline{y} - \underline{a}'_2\|^2.$$

If $\underline{\lambda}$ is a point on the nearest-neighbor decision surface, then the condition of the theorem provides two equations of the form

$$\|\underline{\lambda} - \underline{a}_1\|^2 = \|\underline{\lambda} - \underline{a}_2\|^2$$

$$\text{and } \|\underline{\lambda} - \underline{a}'_1\|^2 = \|\underline{\lambda} - \underline{a}'_2\|^2.$$

We thus have four criteria on $\underline{\lambda}$ and \underline{y} , which are of the form

considered in Lemma 4.2. This provides the contradiction,
and hence the theorem.

Q.E.D.

REFERENCES

1. Batchelor, B. G., "A comparison of the decision surfaces of the nearest-neighbor and potential function classifiers", Information Sciences, Vol. 5, pp. 171-178, 1973.
2. Batchelor, B. G., "Instability of the decision surfaces of the nearest-neighbor and potential function classifiers", Information Sciences, Vol. 5, pp. 179-186, 1973.
3. Batchelor, B. G., Practical approach to pattern classification, Plenum Press, London & New York, 1974.
4. Chahal, J. S., and L. J. White, "Bilinear surface representation of radar scattering for aircraft", Report 2768-10, July 1973, The Ohio State University ElectroScience Laboratory, Department of Electrical Engineering.
5. Chen, C. H., Statistical pattern recognition, Hayden, N. J., 1973.
6. Dasarathy, B., "Some maxmin location and pattern separation problems: theory and algorithms", Ph.D. dissertation, August 1975, The Ohio State University, Department of Computer and Information Science.
7. Duda, R. O., and P. E. Hart, Pattern classification and scene analysis, Wiley, N. Y., 1973.
8. Feller, W., An introduction to probability theory and its applications, Vols. 1 and 2, Third Edition, Wiley, 1968.
9. Fu, K. S., Sequential methods in pattern recognition and machine learning, Academic Press, 1968.
10. Fukunaga, K., Introduction to statistical pattern recognition, Academic Press, 1972.

11. Fukunaga, K. and T. F. Krile, "Calculation of Bayes recognition error for two multivariate Gaussian distributions", IEEE Transactions on Computers, Vol. C-18, No. 3, pp. 220-229, March 1969.
12. Gerald, C. F., Applied numerical analysis, Addison-Wesley, 1973.
13. Hahn, S. F. and G. D. Nelson, "A comparative study of non-parametric methods for pattern recognition", Report SDSU-RSI-72-19, Remote Sensing Institute, South Dakota State University, November 1972.
14. Highleyman, W. H., "The design and analysis of pattern recognition experiments", Bell System Technical Journal, Vol. 41, pp. 723-744, 1962.
15. Hu, M. K., "Visual pattern recognition by moment invariants", IRE Transactions on Information Theory, Vol. IT-8, No. 2, pp. 179-187, February 1962.
16. Kailath, T., "The divergence and Bhattacharyya distance measures in signal selection," IEEE Transactions on Communication Technology, Vol. COM-15, No. 1, pp. 52-60, February 1967.
17. Kanal, L., "Patterns in pattern recognition: 1968-1974", IEEE Transactions on Information Theory, Vol. IT-20, No. 6, pp. 697-722, November 1974.
18. Kanal, L., and B. Chandrasekaran, "On dimensionality and sample size in statistical pattern classification", Pattern Recognition, Vol. 3, pp. 225-234, 1971.
19. Kanal, L., and B. Chandrasekaran, "On linguistic, statistical, and mixed models for pattern recognition", Frontiers of pattern recognition, Ed. S. Watanabe, Academic Press, 1972.
20. Knuth, D. E., The art of computer programming, Vol., 1, Second Edition, Addison-Wesley, 1973.

21. Lachenbruch, P.A., C. Sneeringer, L. T. Revo, "Robustness of the linear and quadratic discriminant function to certain types of non-normality", Communications in Statistics, 1(1), pp. 39-56, 1973.
22. Lin, Y. T., "Classification of objects with complex geometric shape by means of low-frequency electromagnetic response", Ph.D. dissertation, Report 3815-1, August 1974, The Ohio State University, Department of Electrical Engineering.
23. Lin, Y. T. and J. H. Richmond, "E.M. modeling of aircraft at low-frequencies", IEEE Transactions on Antennas and Propagation, Vol. AP-23, No. 1, January 1975.
24. McGhee, R. B., "Automatic recognition of complex three-dimensional objects from optimal images", in Learning systems and intelligent robots, Eds. K. S. Fu and J. T. Tou, Plenum Press, 1974.
25. Naylor, T. H., Computer simulation experiments with models of economic systems, Wiley, 1971.
26. Repjar, A. G., "The linear separability of multiple frequency radar returns with applications to target classification", Ph.D. dissertation, Report 2768-5, November 1970, The Ohio State University, Department of Electrical Engineering.
27. Repjar, A. G., A. A. Ksienski, and L. J. White, "Object identification from multi-frequency radar returns", The Radio and Electronic Engineer, Vol. 45, No. 4, pp. 161-167, April 1975.
28. Roberts, R. A., and C. T. Mullis, "A Bayes sequential test of m hypotheses", IEEE Transactions on Information Theory, Vol. IT-16, No. 1, pp. 91-94, January 1970.
29. Sebestyen, G. S., Decision-making processes in pattern recognition, Macmillan, New York, 1962.

30. Srihari, S. N., and L. J. White, "Pattern recognition techniques for radar aircraft identification", Report 3815-2, October 1975, The Ohio State University ElectroScience Laboratory, Department of Electrical Engineering.
31. Toussaint, G. T., "Bibliography on estimation of misclassification", IEEE Transactions on Information Theory, Vol. IT-20, pp. 472-479, July 1974.
32. Ullmann, J. R., Pattern recognition techniques, Crane-Russak, N. Y., 1973.
33. Wald, A., Sequential analysis, Wiley, 1947.
34. Whalen, A. D., Detection of signals in noise, Academic Press, 1971.
35. White, L. J. and A. A. Ksienski, "Aircraft identification using bilinear surface representation of radar data", Pattern Recognition, Vol. 6, pp. 35-45, 1974.
36. Wong, E. and J. A. Steppe, "Invariant recognition of geometric shapes", in Methodologies of pattern recognition, Ed. S. Watanabe, pp. 535-546, 1969.