

General Disclaimer

One or more of the Following Statements may affect this Document

- This document has been reproduced from the best copy furnished by the organizational source. It is being released in the interest of making available as much information as possible.
- This document may contain data, which exceeds the sheet parameters. It was furnished in this condition by the organizational source and is the best copy available.
- This document may contain tone-on-tone or color graphs, charts and/or pictures, which have been reproduced in black and white.
- This document is paginated as submitted by the original source.
- Portions of this document are not fully legible due to the historical nature of some of the material. However, it is the best reproduction available from the original submission.



DEPARTMENT OF MATHEMATICS
UNIVERSITY OF HOUSTON

HOUSTON, TEXAS

NASA CR-

151823

FEATURE SELECTION FOR BEST MEAN
SQUARE APPROXIMATION OF
CLASS DENSITIES
BY CHARLES PETERS
REPORT #71 JULY 1978

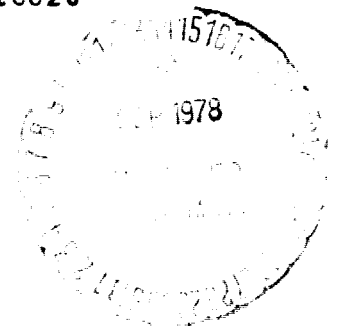
(NASA-CR-151823) FEATURE SELECTION FOR BEST
MEAN SQUARE APPROXIMATION OF CLASS DENSITIES
(Houston Univ.) 12 p HC A02/MF A01 CSCL 12A

N78-30888

Unclas
G3/65 28620

PREPARED FOR
EARTH OBSERVATION DIVISION, JSC
UNDER
CONTRACT NAS-9-15543

HOUSTON, TEXAS 77004



**FEATURE SELECTION FOR BEST MEAN
SQUARE APPROXIMATION OF CLASS DENSITIES**

**Charles Peters*
Department of Mathematics
University of Houston
Houston, Texas**

Report #71

July 1978

*** Author partially supported by NASA/Johnson Space Center under contract NAS-9-15000.**

FEATURE SELECTION FOR BEST MEAN SQUARE APPROXIMATION OF CLASS DENSITIES

ABSTRACT. A criterion for linear feature selection is proposed which is based on mean square approximation of class density functions. It is shown that for the widest possible class of approximants, the criterion reduces to Devijver's Bayesian distance. For linear approximants the criterion is equivalent to well known generalized Fisher criteria.

Pattern recognition

Feature selection

Discriminant analysis

Pattern class separability

Feature Selection for Best Mean Square

Approximation of Class Densities

1. Introduction

The purpose of this note is to describe a general mean square approach to linear feature selection which connects certain generalized Fisher criteria in discriminant analysis with a measure of pattern class separation introduced by Devijver⁽³⁾. The former are typical of those criteria which utilize only low order information about the pattern class distributions, while the latter requires that the class distributions be known, or at least accurately estimated.

Let X denote a random vector in real n -space R^n which arises from one of m pattern classes Π_1, \dots, Π_m having known prior probabilities $\alpha_1, \dots, \alpha_m$, where $\alpha_i > 0$ and $\sum_{i=1}^m \alpha_i = 1$. Let $F_j(x)$ denote the j^{th} class conditional distribution function of X and let $F(x) = \sum_{i=1}^m \alpha_i F_i(x)$ denote the mixture distribution. For a given measurable transformation $T: R^n \rightarrow R^k$ let $G_j(y, T)$ and $G(y, T)$ denote, respectively, the j^{th} class conditional distribution and mixture distribution of the random variable $Y = TX$. We let $f_j(x)$ (resp. $g_j(y, T)$) denote the class conditional densities of X (resp. Y) with respect to their corresponding mixture distributions; i.e.,

$$f_j = \frac{dF_j}{dF} \quad \text{and} \quad g_j(\cdot, T) = \frac{dG_j(\cdot, T)}{dG(\cdot, T)}$$

We will restrict our attention to the set of linear transformations T of rank k , and assume that each pattern class Π_i has a mean μ_i and positive definite covariance matrix Ω_i . Let $\mu = \sum_{i=1}^m \alpha_i \mu_i$ and let

$$S_B = \sum_{i=1}^m \alpha_i (\mu_i - \mu)(\mu_i - \mu)^T$$

$$S_W = \sum_{i=1}^m \alpha_i \Omega_i$$

and

$$S = S_W + S_B$$

denote the between class scatter matrix, the average within class scatter matrix, and the total scatter matrix respectively.

A number of interesting feature selection criteria can be formulated using only the parameters μ , μ_i , Ω_i , S , S_W , S_B ; e.g., the criteria proposed by Kittler and Young⁽⁸⁾, Foley and Sammon⁽⁴⁾, Fukunaga and Koontz⁽⁶⁾, and the discrete analogue of the modified Karhunen-Loeve expansion of Chien and Fu⁽¹⁾. The modified K-L expansion minimizes an entropy function, and also best represents the pattern vector X in an overall least squares sense; however, its value for discrimination has been questioned by several authors (see Kittler⁽⁷⁾). Fukunaga⁽⁵⁾ considers several criteria of the generalized Fisher type, including

$$J_k(T) = \text{tr}(T^T S_W T)^{-1} (T^T S_B T).$$

Thus, according to this criterion, the best $k \times n$ matrix T of rank k is one which maximizes $J_k(T)$. The solution is any T which is row equivalent to a $k \times n$ matrix whose rows are linearly independent principal eigenvectors (i.e., corresponding to the largest eigenvalues) of $S_W^{-1} S_B$. We also consider a modification

$$J_k'(T) = \text{tr} (T^T S T)^{-1} (T^T S_B T)$$

which admits the same maximizing T .

The Bayesian distance corresponding to the pattern classes Π_1, \dots, Π_m , as defined by Devijver⁽³⁾, is

$$\begin{aligned} B_n &= \sum_{i=1}^m \alpha_i^2 E [f_i(X)^2] \\ &= \sum_{i=1}^m \alpha_i^2 \int_{R^n} f_i(x)^2 dF(x). \end{aligned}$$

Its transformed value is

$$B_k(T) = \sum_{i=1}^m \alpha_i^2 \int_{R^k} g_i(y, T)^2 dG(y, T).$$

Devijver proves a number of interesting inequalities relating B_n to the Bayes probability of misclassification, the Bhattacharyya coefficient, and other measures of class separation. In addition, he notes that Cover and Hart⁽²⁾ have shown that $1 - B_n$ is the asymptotic error rate of the nearest neighbor classifier.

2. Mean Square Optimality of Bayesian Distance

For a given $k \times n$ matrix T of rank k , let $L_2(T)$ denote the set of all measurable functions $\varphi : R^k \rightarrow R^1$ such that $\int_{R^k} \varphi(y)^2 dG(y, T) < \infty$ and let C_T be a given closed linear subspace of $L_2(T)$. Our general approach to linear feature selection is to choose that T , if possible, which minimizes

$$R(T) = \sum_{i=1}^m \beta_i \min_{\varphi_i \in C_T} \int_{R^n} [\varphi_i(Tx) - f_i(x)]^2 dF(x),$$

where the β_i are positive weights. That is, we attempt to find a T which produces a set of approximations $\varphi_i(Tx)$ to the class densities $f_i(x)$ which is best in an overall mean square sense. Given such approximations we may classify observations of X according to the pseudo-Bayes rule: decide that X is from class Π_i if $\alpha_i \varphi_i(Tx) > \alpha_j \varphi_j(Tx)$ for each $j \neq i$. Since we are interested in classification accuracy, it seems appropriate to choose weights β_i which reflect the relative importance of the classes in the mixture distribution; e.g., $\beta_i = \alpha_i$ for all i or $\beta_i = \alpha_i^2$ for all i . For the remainder of this section we choose $\beta_i = \alpha_i^2$ and $C_T = L_2(T)$.

Proposition 1: For $\beta_i = \alpha_i^2$, $i = 1, \dots, m$ and $C_T = L_2(T)$ for each T ,

$$R(T) = B_n - B_k(T).$$

Proof: Observe that $g_i(y, T) \in L_2(T)$, since it is bounded by α_i^{-1} . Moreover, for each $\varphi \in L_2(T)$,

$$\begin{aligned} & \int_{R^n} \varphi(Tx) [g_i(Tx, T) - f_i(x)] dF(x) \\ &= \int_{R^n} \varphi(Tx) g_i(Tx, T) dF(x) - \int_{R^n} \varphi(Tx) dF_i(x) \\ &= \int_{R^k} \varphi(y) g_i(y, T) dG(y, T) - \int_{R^k} \varphi(y) dG_i(y, T) \\ &= 0. \end{aligned}$$

Therefore,

$$\begin{aligned}
 \min_{\varphi \in L_2(T)} \int_{\mathbb{R}^n} [\varphi(Tx) - f_1(x)]^2 dF(x) \\
 &= \int_{\mathbb{R}^n} [g_1(Tx, T) - f_1(x)]^2 dF(x) \\
 &= \int_{\mathbb{R}^n} f_1(x)^2 dF(x) - \int_{\mathbb{R}^n} g_1(y, T)^2 dG(y, T).
 \end{aligned}$$

The assertion of the proposition follows on multiplying by α_1^2 and summing over i .

We may summarize by saying that if there exists a $k \times n$ matrix T_0 of rank k which maximizes $B_k(T)$, then the functions $g_1(T_0x, T_0), \dots, g_m(T_0x, T_0)$ constitute the best mean square approximation to the class densities $f_1(x), \dots, f_m(x)$ attainable through a linear compression of the data into k dimensions. Since $B_k(QT) = B_k(T)$ for each nonsingular $k \times k$ matrix Q and each $k \times n$ matrix T of rank k , $B_k(T)$ has a maximum if and only if it has a maximum on the compact set $\{T \mid TT^T = I_{k \times k}\}$. In particular, if $B_k(T)$ is continuous, it has a maximum.

3. Best Linear Approximation of f Class Densities

In this section we let C_T be the set of functions $\varphi(y) = w + b^T y$, where w is a real number and $b \in \mathbb{R}^k$. For simplicity, we use the notation

$$\varphi(y) = a^T v(y), \text{ where } a = \begin{pmatrix} w \\ b \end{pmatrix} \in \mathbb{R}^{k+1} \text{ and } v(y) = \begin{pmatrix} 1 \\ y \end{pmatrix} \in \mathbb{R}^{k+1}. \text{ For given } T,$$

$$a_1 = \begin{pmatrix} w \\ b \end{pmatrix} \text{ minimizes}$$

$$\begin{aligned}
& \int_{\mathbb{R}^n} [a^T v(Tx) - f_1(x)]^2 dF(x) \\
&= a^T \left[\int_{\mathbb{R}^n} v(Tx) v(Tx)^T dF(x) \right] a \\
&\quad - 2a^T \int_{\mathbb{R}^n} v(Tx) dF_1(x) + \int_{\mathbb{R}^n} f_1(x)^2 dF(x) \\
&= a^T \left(\begin{array}{c|c} 1 & \mu^T T^T \\ \hline T\mu & TWT^T \end{array} \right) a - 2(1 \mid \mu^T T^T) a \\
&\quad + \int_{\mathbb{R}^n} f_1(x)^2 dF(x)
\end{aligned}$$

if and only if

$$\left(\begin{array}{c|c} 1 & \mu^T T^T \\ \hline T\mu & TWT^T \end{array} \right) \begin{pmatrix} w_1 \\ b_1 \end{pmatrix} = \begin{pmatrix} 1 \\ T\mu \end{pmatrix}$$

where $W = E[XX^T] = S + \mu\mu^T$. Solving this system gives

$$w_1 = 1 - \mu^T T^T (TST^T)^{-1} T(\mu_1 - \mu)$$

and

$$b_1 = (TST^T)^{-1} T(\mu_1 - \mu).$$

The corresponding squared error of approximation is

$$-1 - (\mu_1 - \mu)^T T^T (TST^T)^{-1} T(\mu_1 - \mu) + \int_{\mathbb{R}^n} f_1(x)^2 dF(x).$$

Therefore, the criterion to be minimized is

$$R(T) = - \sum_{i=1}^m \beta_i (\mu_i - \mu)^T T^T (TST^T)^{-1} T (\mu_i - \mu)$$

+ terms independent of T.

That is, we want to maximize

$$\hat{R}(T) = \text{trace } (TST^T)^{-1} \hat{TS}_B T^T,$$

where

$$\hat{S}_B = \sum_{i=1}^m \beta_i (\mu_i - \mu) (\mu_i - \mu)^T.$$

The solution is $T = QT_0$, where T_0 is a $k \times n$ matrix whose rows are linearly independent principal eigenvalues of, $S^{-1} \hat{S}_B$ and Q is an arbitrary nonsingular $k \times k$ matrix. In particular, for $\beta_i = \alpha_i$ we obtain the same solution given by Fukunaga's criterion,

$$\text{trace } (TS_w T^T)^{-1} (TS_B T^T).$$

4. Concluding Remarks

An equivalent set of criteria for feature selection are expressions such as

$$\bar{R}(T) = \sum_{i=1}^m \bar{\beta}_i \min_{\varphi \in C_T} \int_{R^n} [\varphi(Tx) - \alpha_i f_i(x)]^2 dF(x)$$

in which the posterior probabilities $\alpha_i f_i(x)$ of the classes are approximated. If each $\bar{\beta}_i$ is chosen to be 1 $\bar{R}(T)$ is the same as $R(T)$ with $\beta_i = \alpha_i^2$. This, together with Proposition 1 and the relationship between Bayesian distance and the probability of error, seems to indicate that the choice $\beta_i = \alpha_i^2$ is a good one.

In some cases it may be numerically feasible to use $B_k(T)$ as a feature selection criterion when assumptions about the parametric form of the class distributions are made. For example, if each class distribution F_i is multivariate normal, then $B_k(T)$ reduces to an expression which is continuously differentiable in T and which, moreover, can be approximated by sample averages over an unlabeled sample from the mixture distribution. Thus descent algorithms might be successfully employed in maximizing $B_k(T)$.

Finally, we remark that there is no reason in principle why $B_k(T)$ cannot be regarded as a criterion for nonlinear feature extraction. Indeed, Proposition 1 remains true when T is any measurable transformation from R^n onto R^k .

5. Acknowledgements

The author was partially supported by NASA Johnson Space Center under contract NAS-9-15000.

REFERENCES

1. Y. T. Chien and K. S. Fu. On the generalized Karhunen-Loeve expansion. IEEE Trans. IT-15, 518-520, (1967).
2. T. M. Cover and P. E. Hart. Nearest neighbor pattern classification. IEEE Trans. IT-13, 21-27, (1967).
3. P. A. Devijver. On a new class of bounds on Bayes risk in multihypothesis pattern recognition. IEEE Trans. C-23, 70-80, (1974).
4. D. H. Foley and J. W. Sammon, Jr. An optimal set of discriminant vectors. IEEE Trans. C-24, 281-289, (1975).
5. K. Fukunaga. Introduction to Statistical Pattern Recognition, Academic Press, New York, (1972).
6. K. Fukunaga and W. L. G. Koontz. Application of the Karhunen-Loeve expansion to feature selection and ordering. IEEE Trans. C-19, 311-318 (1970).
7. J. Kittler. Mathematical methods of feature selection in pattern recognition. Int. J. Man-Machine Studies, 7, 609-637, (1975).
8. J. Kittler and P. C. Young. A new approach to feature selection based on the Karhunen-Loeve expansion, Pattern Recognition, Vol. 5, 335-352, (1973).