

## PENALIZED DISCRIMINANT ANALYSIS

BY TREVOR HASTIE, ANDREAS BUJA AND ROBERT TIBSHIRANI<sup>1</sup>

*Stanford University, AT & T Bell Laboratories  
and University of Toronto*

Fisher's linear discriminant analysis (LDA) is a popular data-analytic tool for studying the relationship between a set of predictors and a categorical response. In this paper we describe a penalized version of LDA. It is designed for situations in which there are many highly correlated predictors, such as those obtained by discretizing a function, or the grey-scale values of the pixels in a series of images. In cases such as these it is natural, efficient and sometimes essential to impose a spatial smoothness constraint on the coefficients, both for improved prediction performance and interpretability. We cast the classification problem into a regression framework via optimal scoring. Using this, our proposal facilitates the use of any penalized regression technique in the classification setting. The technique is illustrated with examples in speech recognition and handwritten character recognition.

**1. Introduction.** Linear discriminant analysis (LDA) is a standard tool both for classification and dimension reduction. Here is roughly how LDA works:

1. Classification is based on the Mahalanobis distances  $(x - m^j)\Sigma_W^{-1}(x - m^j)$  to class centroids  $m_j$ , where  $\Sigma_W$  is the pooled within-class covariance of the predictors  $x$ .
2. Data reduction entails a sequence of unit-variance, linear-discriminant variables  $\beta_k^T x$ , chosen to maximize successively  $\beta_k^T \Sigma_{\text{Bet}} \beta_k$ , with  $\Sigma_{\text{Bet}}$  the between-class covariance matrix. These discriminant variables represent a subspace for which the class centroids are spread out as much as possible.

Linear discriminant analysis enjoys a number of favorable properties, such as reasonable robustness to nonnormality and even to mildly different class covariances [see the many references in Seber (1984), Chapter 6]. On the down side, there are two deficiencies which apply under opposite circumstances:

1. LDA is too flexible in situations with large numbers (e.g., hundreds) of highly correlated predictor variables;
2. it is too rigid in situations where the class boundaries in predictor space are complex and nonlinear.

In the first case, LDA overfits, in the second case, LDA underfits the data.

---

Received August 1992; revised May 1994.

<sup>1</sup>Supported by the Natural Sciences and Engineering Research Council of Canada.

AMS 1991 *subject classifications*. Primary 62H30; secondary 62G07.

*Key words and phrases*. Signal and image classification, discrimination, regularization.

We show that both problems can be overcome by modifications of LDA that effectively regularize a large, nearly or fully degenerate within-class covariance matrix  $\Sigma_w$ . This paper focuses on the first situation where LDA has problems: large numbers of correlated predictor variables. In a companion paper [Hastie, Tibshirani and Buja (1994)], we address the other problem, interestingly by using results of this paper.

We call the technique discussed here *penalized discriminant analysis*, or PDA for short. We rely on the relationship between linear discriminant analysis and canonical correlation analysis, or more precisely its asymmetric cousin, here called optimal scoring. If viewed appropriately, optimal scoring contains *linear regression* as a building block. This building block lends itself to generalization by replacing linear least squares regression with other types of regression. In the present paper, we use *penalized least squares regression* to overcome the problems of high-dimensional (e.g., more than 200) correlated predictors, while in the companion paper we use nonparametric adaptive regression methods to solve the problem of complex class boundaries when the predictors are relatively low-dimensional (probably no more than 30).

Along with the transfer of regression methods to discriminant analysis goes the transfer of related notions. For example, we can now use degrees of freedom in classification problems; these arise naturally when selecting or comparing smoothing parameters in fixed-bandwidth nonadaptive regression methods for optimal scoring.

In this paper we study two examples with large numbers of correlated predictor variables:

1. In speech recognition, one is interested in classifying short speech frames into one of several phoneme classes, based on the log-periodogram of the frame (for example); a typical log-periodogram estimate forms a predictor vector of dimension 256.
2. In handwritten character recognition, one wishes to classify small images of characters and digits into the obvious classes; a typical size-normalized image might be represented by 16-by-16 gray-scale pixels, and again form a predictor of dimension 256.

The data in both of these examples arise from the discretization of analog signals. It is obvious that an empirical (within-class) covariance matrix  $\Sigma_w$  of size 256 by 256 will have unfavorable statistical properties for almost any realistic size of the training sample, and straightforward Mahalanobis distances calculated in LDA will suffer or become useless as a result. Although we used a discretization resolution of 256 in both cases, one could use a higher resolution and the problems would be worse. Some form of “borrowing strength” or regularization appears to be needed and can be formally motivated in the context of a population model [see Section 4 as well as Leugrants, Moyeed and Silverman (1993)].

Even if we have sufficient data, the estimates are likely to be spatially rough. To see this, consider the usual linear discriminant functions with

coefficient vectors  $\Sigma_W^{-1}m^j$ , where  $m^j$  is the mean vector for the  $j$ th class. If adjacent predictor (log-spectral or gray-scale) values have strong positive correlations due to locally smooth behavior, then the spectral decomposition of  $\Sigma_W$  will favor low-frequency functions (large eigenvalues for smooth eigenvectors, small eigenvalues for rough eigenvectors);  $\Sigma_W^{-1}$  will have the inverse structure, and hence  $\Sigma_W^{-1}m^j$  will emphasize the rough components of  $m^j$ . This phenomenon has two undesirable consequences: (1) rough coefficient contours that lack smoothness on the index domain are not interpretable; (2) misclassification rates deteriorate since jagged coefficient contours indicate reliance on irrelevant local contrasts that are easily thrown off if the behavior of a test predictor vector deviates slightly from those found in the training sample.

Similar problems have been recognized by several authors [Di Pillo (1976, 1979), Campbell (1980) and Friedman (1989)] who introduce ridge-type regularizations of LDA. The obvious idea is to stabilize the (within-class) covariance matrix by adding a diagonal matrix, often a small multiple of the identity. The relation of this approach to ridge regression goes beyond a formal analogy due to the above-mentioned relation between LDA and optimal scoring: subjecting the regression building block of optimal scoring to a ridge-type modification gives essentially the regularized version of LDA proposed by these authors. Also equivalent is Vinod's (1976) ridge-regularization of canonical correlation analysis when applied to a discriminant context [see also Vinod and Ullah (1981)]. Vinod's aim was to counter the detrimental effects of collinear sets of variables on canonical correlation analysis.

It appears, though, that the examples of speech and character recognition call for a different type of regularization. We often wish to interpret the coefficient arrays of discriminant variables, in particular in the two cases above. These coefficient arrays can be interpreted as discretized curves and images since they, too, are indexed by frequency and pixel location, respectively. In order to achieve spatially smooth behavior, the usual ridge method does not seem entirely appropriate since its bias toward an overall mean ignores the spatial structure of the index domain. It may therefore be more plausible to bias the discriminant coefficients toward smoothness as a function of the frequency or the pixel location, for example, through regularization along the lines of smoothing spline models. Smoothing splines enforce smoothness by penalizing local contrasts such as second-order differences, that is, the coefficient array of a discriminant variable pays a price for locally rough behavior as a function of frequency or pixel location.

To summarize, two distinct motivations for regularization have emerged:

1. When the number-of-variables to sample-size ratio is too high, we cannot reliably estimate a covariance matrix. Since the large number of variables arise from discretizing an analog signal, natural methods of regularization can be found.
2. Even if the sample size were sufficient to estimate the structured covariance matrix, coefficients of spatially smooth variables tend to be spatially

rough. Since we hope to interpret these coefficients, we would prefer smoother versions, especially if they do not compromise the fit.

Our proposal replaces  $\Sigma_W$  by a regularized version  $\Sigma_W + \lambda\Omega$ , where  $\Omega$  is a “roughness”-type penalty matrix; the LDA analysis then proceeds as usual. Figures 1 and 2 show the results for the two examples. In each we see the raw LDA coefficients displayed as curves and images, respectively, as well as the corresponding regularized versions. In both cases the regularized versions improved classification performance on test data significantly. Both these examples are treated in detail in Sections 5 and 6.

The possibility of using other than ridge-type penalties was mentioned in passing by Friedman [(1989), end of Section 8]. He proposes to penalize local averages to enforce smoothness, which is distinct from our proposal since the splines we use penalize local differences. We do not know whether Friedman’s proposal achieves what it sets out to do. The thrust of his paper is to combine

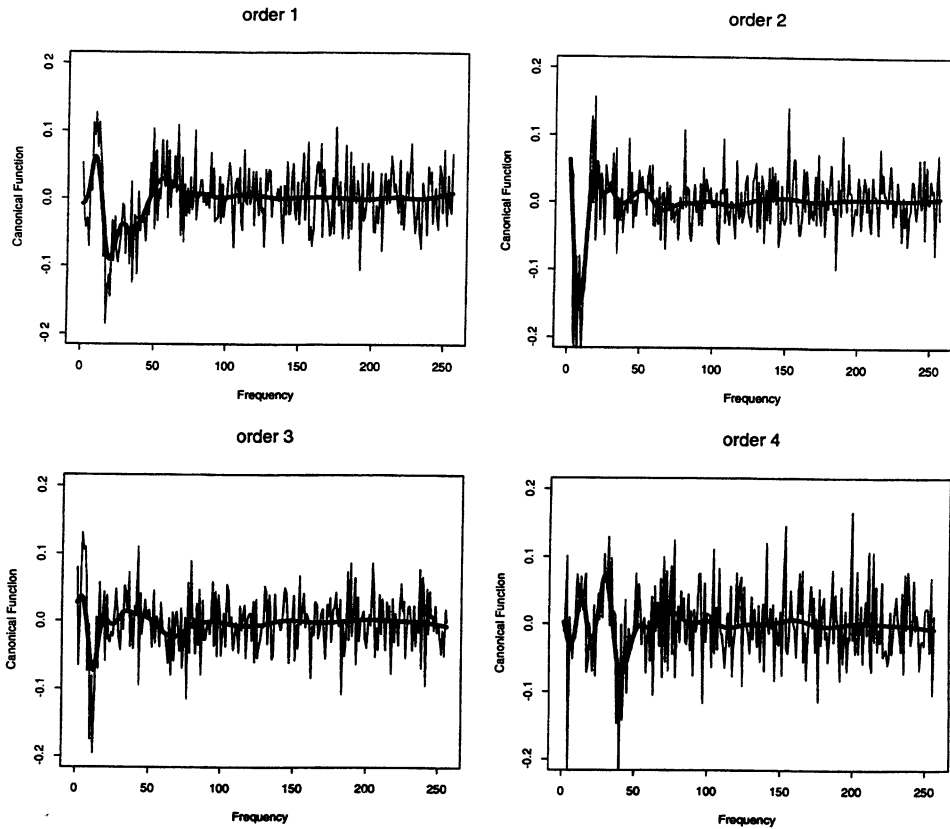


FIG. 1. *The jagged functions are the four discriminant coefficient functions for separating the five phonemes in the sampled utterances in the example of Section 5. The superimposed smooth functions are the regularized discriminant coefficient functions, using 30 df worth of smoothing.*

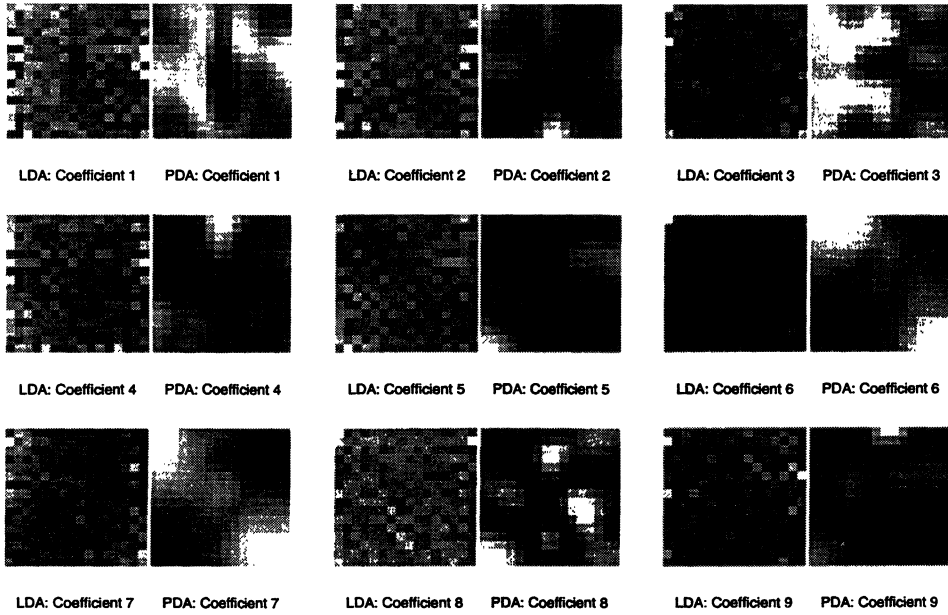


FIG. 2. The images appear in pairs and represent the nine discriminant coefficient functions for the digit recognition problem. The left member of each pair is the LDA coefficient, while the right member is the PDA coefficient, regularized to enforce spatial smoothness.

ridge-shrinkage, LDA and quadratic discriminant analysis (QDA) in a single framework; that is, he considers problems of lower dimensionality than we do since QDA estimates within-class covariances individually for each class, thus further compounding the problem caused by high dimensionality.

A different and quite complex approach is taken by Kiiveri (1992), whose motivation also arose from the analysis of high-dimensional spectral data. He assumes a factor-analysis-growth-curve model for the within-class covariance:  $\Sigma_W = \Lambda\Phi\Lambda^T + \Psi$ , where the first term is structural and of low rank, and  $\Psi$  is typically diagonal, stemming from uncorrelated errors. Based on normal assumptions, the components are estimated with an EM algorithm. As an aside, Kiiveri [(1992), end of Section 2.2] mentions the possibility of using a stationary rather than diagonal form for  $\Psi$ . His  $\Psi$  is reminiscent of a penalty, but its actual role is that of a factor analytic “uniqueness.” He seems exclusively interested in interpretability of canonical variates since he does not mention misclassification rates.

Related to the present work is Leurgans, Moyeed and Silverman (1993), who study canonical-correlation analysis of pairs of discretized analog signals. They also use a penalty approach and they provide some asymptotic theory. More details are given in Section 4.

We should also mention an approach that seems obvious but does not work: motivated by the fact that log-spectra are often smoothed before

further processing, one might argue that LDA should be applied to smoothed log-spectra and, by generalization, to smoothed images. The problem with this approach is that smoothing adds to already existing correlations among neighboring log-spectral and gray-scale values, thus causing the canonical variate arrays to be even more jagged. The within-class covariance of smoothed data will be even more degenerate than that of the raw data, thus calling for even more regularization. Thus we see that regularization at the level of the within-class covariance is indispensable. Presmoothing then turns out to be superfluous, adding to the computational cost of PDA with zero benefit in classification performance.

Another approach to regularization is *filtering*, popular in the signal and speech processing literature. Here the spatial structure of the predictors is acknowledged by approximating the data by its projection onto a low-dimensional, spatially smooth basis. The predictors are thus replaced by their basis coefficients. A drawback of this approach is that one has to supply a suitable basis, which tends to involve more subjectivity than choosing a smoothing penalty. In our examples the filtering approach produced results similar but slightly inferior to those achieved by regularization.

**2. Notation and conventions.** Consider a discrimination problem with  $J$  classes and  $N$  training samples. The training samples consist of observed predictor vectors and their known class memberships. These memberships will be represented by a categorical response variable  $G$  with  $J$  levels; the individual responses, by  $g_i$ . It is sometimes convenient to code the  $N$  responses in terms of an indicator matrix  $Y$  of size  $N \times J$ , with columns that correspond to the dummy-variable codings of the  $J$  classes. We denote by  $h_i$  the  $m$ -vector of predictor values (input features) for the  $i$ th training sample. In the digit-recognition example,  $h_i$  is a 256-vector of gray-scale pixel measurements for the  $i$ th training image. In the context of flexible discriminant analysis [Hastie, Tibshirani and Buja (1994)], the  $h_i$  will be expanded bases (splines, polynomials, etc.) of observed predictor variables  $x_i$ , that is,  $h_i = h(x_i)$ . The vectors  $h_i$  for all the training observations will be stored as rows of the  $N \times m$  matrix  $H$ , which we call the predictor matrix. We will assume that the columns of  $H$  are centered (i.e., orthogonal to the constant 1-vector). If this assumption is not satisfied, the constant 1-vector should be an element of the column space of  $H$ , in which case the eigenproblems encountered below will exhibit a trivial eigenvalue which is easily weeded out.

We assume that some penalty matrix  $\Omega$  of size  $m \times m$  is given, for penalizing the coefficients of  $h$ . It should be symmetric and nonnegative definite. If a smoothing parameter  $\lambda$  is part of the penalty, we assume that it is absorbed in  $\Omega$ . These penalty matrices arise naturally in the context of the examples. If the data are log-spectra or images, we define  $\Omega$  in such a way as to force nearby components of  $\beta_k$  to be similar. If  $h_i = h(x_i)$  represent a basis expansion of the actual input vectors  $x_i$ , a penalty matrix may be chosen so that the compositions  $h(x_i)^T \beta_k$  are smooth as functions of  $x$ . See Sections 5, 6 and 7.

**3. The equivalence of penalized linear discriminant analysis, canonical correlation analysis and optimal scoring.** It is known that discriminant variates are up to scale factors the same as the so-called canonical variates which result from an associated canonical correlation analysis (CCA), and often the latter term is used interchangeably with discriminant variates. Somewhat lesser known is that an asymmetric version of canonical correlation analysis, here called *optimal scoring* and abbreviated OS, also yields a set of dimensions which coincide up to scalars with those of LDA and CCA.

In the following sections we introduce OS, CCA and LDA with penalization built in. We then show that the three are equivalent just as they are without penalization. This has the important consequence that we can simply use our tools for penalized and nonparametric regression to perform penalized and nonparametric discriminant analysis.

We also discuss the dimension reduction aspect of linear discriminant analysis. Dimension reduction means reexpressing the data in fewer variables while minimizing the loss of essential information for the problem at hand. Such reduction can actually be beneficial when the “lost dimensions” show only spurious or weak structure. The reduced dimensions resulting from LDA are variously called discriminant variates, discriminant coordinates or sometimes crimcoords for short.

Each of the three problems (OS, CCA and LDA) to be defined below has an associated criterion and constraints under which the criterion is to be optimized or made stationary.

**3.1. Penalized optimal scoring.** The point of optimal scoring is to turn categorical variables into quantitative ones by assigning scores to classes (groups, categories). In our notation above,  $\theta(G)$  assigns a real number (say,  $\theta_j$ ) to the  $j$ th level of  $G$ . Given a  $J$ -vector of such scores  $\theta$  for the  $J$  classes, the  $N$ -vector  $Y\theta$  represents a vector of scored training data which one may try to regress onto the predictor matrix  $H$ . (With a slight abuse of notation, we use  $\theta$  to represent both the function or the vector of  $J$  real numbers that represent the function.) The simultaneous estimation of scores and regressions constitutes the optimal scoring problem:

**DEFINITION 1.** The penalized optimal scoring problem is defined by the criterion

$$(1) \quad \text{ASR}(\theta, \beta) = N^{-1} \left( \sum_{i=1}^N [\theta(g_i) - h(x_i)^T \beta]^2 + \beta^T \Omega \beta \right)$$

$$(2) \quad = N^{-1} (\|Y\theta - H\beta\|^2 + \beta^T \Omega \beta),$$

which is to be minimized (made stationary) under the constraint  $N^{-1} \|Y\theta\|^2 = 1$ .

Although Definition 1 is stated in terms of a single solution  $(\theta, \beta)$ , implicit is a sequence of solutions  $(\theta_k, \beta_k)$  with orthogonality defined by the implied

inner product  $N^{-1}\langle Y\theta_k, Y\theta_l \rangle = \delta_{kl}$ . We use this compact notation in the subsequent definitions as well, and thus we avoid a more cumbersome notation involving traces of matrices.

Optimal scoring (or scaling) is common in correspondence analysis [e.g., Lebart, Morineau and Warwick (1984); see Chapter 3 for a comparison with discriminant analysis] and the psychometric literature [e.g., Gifi (1981, 1990), in particular Chapter 7.2 for discriminant analysis; and the series of papers by de Leeuw, Takane, Young, in various permutations (1976, 1978, 1979)].

It is useful to interpret criterion (2) as a quadratic form in the combined  $\theta$ - $\beta$  vector:

$$(3) \quad \text{ASR}(\theta, \beta) = \theta^T \Sigma_{11} \theta - 2\theta^T \Sigma_{12} \beta + \beta^T \Sigma_{22} \beta,$$

where the obvious definitions are as follows:

1.  $\Sigma_{11} = N^{-1}Y^T Y$  is a diagonal matrix with the class proportions  $p_i = N_i/N$  in the diagonal;
2.  $\Sigma_{22} = N^{-1}(H^T H + \Omega)$  is the penalized covariance matrix of the predictors;
3.  $\Sigma_{12} = N^{-1}Y^T H$ ;  $\Sigma_{21} = \Sigma_{12}^T$ .

If we assume that all classes are nonempty ( $N_j > 0$ ),  $\Sigma_{11}$  is invertible. As far as  $\Sigma_{22}$  is concerned, we must assume that the penalty  $\Omega$  is chosen intelligently so as to prevent exact or near collinearities among the predictors, that is,  $\Sigma_{22}$  should be invertible, too.

For a given score vector  $\theta$ , the minimizing  $\beta$  for the OS problem is the penalized least squares estimate:

$$(4) \quad \beta_{os} = (H^T H + \Omega)^{-1} H^T Y \theta = \Sigma_{22}^{-1} \Sigma_{21} \theta,$$

and the partially minimized criterion becomes:

$$(5) \quad \min_{\beta} \text{ASR}(\theta, \beta) = 1 - N^{-1} \theta^T Y^T S(\Omega) Y \theta = 1 - \theta^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \theta,$$

where  $S(\Omega) = H(H^T H + \Omega)^{-1} H^T$  denotes the ‘‘hat’’ or ‘‘smoother’’ matrix of  $H$  regularized by  $\Omega$ . This form of the OS problem is computationally useful because it shows that one can run a penalized multiresponse regression  $Y$  onto  $H$  once:  $\hat{Y} = S(\Omega)Y$ , and then eigenanalyze  $Y^T \hat{Y}$  (with respect to the metric  $\Sigma_{11}$ ) to obtain the stationary  $\theta$  vectors. In Section 7 we discuss various forms for  $S(\Omega)$ . The intent of the following sections is to translate the results of OS analysis into a LDA form.

### 3.2. Penalized canonical correlation analysis.

**DEFINITION 2.** The penalized canonical correlation problem is defined by the criterion

$$\text{COR}(\theta, \beta) = \theta^T \Sigma_{12} \beta,$$

which is to be maximized (made stationary) under the constraints

$$\theta^T \Sigma_{11} \theta = 1 \quad \text{and} \quad \beta^T \Sigma_{22} \beta = 1.$$



Formally, this looks like the usual canonical correlation problem as applied to linear discrimination, except for the penalization built into  $\Sigma_{22}$ .

The criteria of optimal scoring and canonical correlation analysis are related to each other by (3): under the CCA constraints,  $ASR = 2 - COR$ , which shows that the two problems differ only in the additional constraint on  $\beta$  which is missing in OS. For a given set of scores  $\theta$ , the maximizing  $\beta$  for the CCA problem is, up to a scalar, the same as the minimizing  $\beta$  (4) for the OS problem:

$$(6) \quad \beta_{CCA} = \beta_{OS} / \sqrt{\beta_{OS}^T \Sigma_{22} \beta_{OS}},$$

which entails

$$(7) \quad \max_{\beta^T \Sigma_{22} \beta = 1} COR(\theta, \beta) = (\theta^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \theta)^{1/2}.$$

A comparison with (5) shows that the OS and CCA problems produce identical stationary score vectors  $\theta$ .

**3.3. Penalized linear discriminant analysis.** Linear discriminant analysis finds linear combinations of the predictors that maximize the between-class variance (of the class means), relative to the within-class variance. To introduce the *penalized* linear discriminant problem, we need the customary decomposition of the *total covariance*  $\Sigma_{22}$  into *between-class covariance* and *within-class covariance* (or respective cross-products, if the columns of the expanded predictor matrix  $H$  are not centered). The between-class covariance is the covariance of  $H$  regressed onto  $Y$ , or, equivalently, the class-weighted covariance of the class means. The within-class covariance is what is left, and is a pooled estimate of the common covariance matrix for the  $J$  classes. Let  $P_Y = Y(Y^T Y)^{-1} Y^T$  be the projector onto  $Y$ -column space:

1.  $M = \Sigma_{11}^{-1} \Sigma_{12}$  is a  $J \times m$  matrix whose rows are the class means  $m^j = \text{ave}(h_i; i \in \text{Class } j)$ :  $M = (m^1, \dots, m^J)^T$ ;
2.  $\Sigma_{\text{Bet}} = N^{-1} (P_Y H)^T (P_Y H) = \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} = M^T \Sigma_{11} M$ ;
3.  $\Sigma_W = N^{-1} [(I - P_Y) H]^T [(I - P_Y) H + \Omega] = \Sigma_{22} - \Sigma_{\text{Bet}}$ .

Accordingly, penalization affects only the within-class covariance.

**DEFINITION 3.** The criterion of the penalized linear discriminant problem is the (unpenalized) between-class variance,

$$BVAR(\beta) = \beta^T \Sigma_{\text{Bet}} \beta,$$

which is to be maximized (made stationary) under a constraint on the penalized within-class variance,

$$WVAR(\beta) = \beta^T \Sigma_W \beta = 1.$$

The equivalence of the CCA and the LDA problem is standard [e.g., 11.5.4 of Mardia, Kent and Bibby (1979)], but it will follow from the following section as well. The interpretation of this particular form of penalization is

more easily made in the context of a particular example. Highly correlated predictors, such as in the case of discretized spectra, can meet the (unpenalized) normalization constraint by using negatively correlated coefficients, since the constraint is in terms of the variance of the derived variable. Viewed as a function, the negatively correlated coefficients will appear wiggly; the appropriate penalty in this case would limit the spatial roughness of these coefficients.

3.4. *Translation of optimal scoring dimensions into discriminant coordinates.* It is convenient to use CCA as a link between OS and LDA; CCA is a generalized singular value problem for  $\Sigma_{12}$  with regard to the metrics given by  $\Sigma_{11}$  and  $\Sigma_{22}$ . The associated singular value decomposition, essentially a collection of stationary solutions of the CCA problem, takes on this form:

$$(8) \quad \Sigma_{11}^{-1} \Sigma_{12} \Sigma_{22}^{-1} = \Theta D_{\alpha} B^T,$$

$$(9) \quad \Theta^T \Sigma_{11} \Theta = I_L,$$

$$(10) \quad B^T \Sigma_{22} B = I_L,$$

where  $L = \min(J, m)$ ,  $\Theta$  is a  $J \times L$  matrix whose columns  $\theta_k$  are left-stationary vectors,  $B$  is an  $m \times L$  matrix whose columns  $\beta_k$  are right-stationary vectors and  $D_{\alpha}$  is a diagonal matrix of size  $L \times L$  with nonnegative diagonal elements  $\alpha_k$  sorted in descending order.

The usual proof of (8) is by coordinate transformations  $\theta = \Sigma_{11}^{-1/2} \theta'$  and  $\beta = \Sigma_{22}^{-1/2} \beta'$ , which bring (9) and (10) to a Euclidean form, so the simple (nongeneralized) SVD can be applied to  $\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1/2}$  (i.e.,  $\Sigma_{12}$  expressed in the new coordinates). A simple SVD of the form  $A = UDV^T$  entails the trivial consequences  $AV = UD$ ,  $A^T U = VD$ ,  $U^T AV = D$ ,  $V^T A^T AV = D^2$ ,  $U^T AA^T U = D^2$ ; these are translated to the generalized SVD as follows:

$$(11) \quad \Sigma_{11}^{-1} \Sigma_{12} B = \Theta D_{\alpha}$$

$$(12) \quad \Sigma_{22}^{-1} \Sigma_{21} \Theta = B D_{\alpha},$$

$$(13) \quad \Theta^T \Sigma_{12} B = D_{\alpha},$$

$$(14) \quad \Theta^T \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Theta = D_{\alpha}^2,$$

$$(15) \quad B^T \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} B = D_{\alpha}^2.$$

In particular, (13) implies  $\text{COR}(\theta_k, \beta_k) = \alpha_k$ . [Here  $D_{\alpha}$  denotes a diagonal matrix with elements  $\alpha_k$ ,  $D_{\alpha}^2$  with diagonal elements  $\alpha_k^2$ ].

According to (5) and (7) the stationary  $\theta$  vectors of OS and CCA are the same, while the  $\beta$  vectors of OS and CCA are related according to (4) and (12) by

$$(16) \quad B_{\text{OS}} = B D_{\alpha},$$

$B_{\text{OS}}$  being a matrix of OS-stationary column vectors  $\beta_{\text{OS}, k}$ . From (5) and (7) it follows that  $\text{ASR}(\theta_k, \beta_k) = 1 - \alpha_k^2$ .

To link CCA and LDA, we translate (15) directly into

$$(17) \quad B^T \Sigma_{\text{Bet}} B = D_{\alpha^2},$$

$$(18) \quad B^T \Sigma_W B = I_L - D_{\alpha^2} = D_{1-\alpha^2}.$$

This shows that  $B$  diagonalizes both  $\Sigma_{\text{Bet}}$  and  $\Sigma_W$ . If we define

$$(19) \quad B_{\text{LDA}} = B D_{(1-\alpha^2)^{-1/2}},$$

we get a matrix whose columns  $\beta_{\text{LDA},k}$  are stationary solutions of the LDA problem:

$$(20) \quad B_{\text{LDA}}^T \Sigma_W B_{\text{LDA}} = I_L, \quad B_{\text{LDA}}^T \Sigma_{\text{Bet}} B_{\text{LDA}} = D_{\alpha^2/(1-\alpha^2)}.$$

Finally, the relation between the LDA and OS solutions is

$$(21) \quad B_{\text{LDA}} = B_{\text{OS}} D_{[\alpha^2(1-\alpha^2)]^{-1/2}}.$$

**3.5. Graphical projections.** The strongest discriminant coordinates are useful for graphical examination of both training and test data, the former visually to assess overlap of groups in predictor space, the latter to judge the strength of evidence for assigning test data to specific classes. The elements that need to be plotted are projections of the following:

1. the class centers  $m^j$ ;
2. the predictor vectors  $h_i$  of the training data with an indication of their class memberships;
3. the predictor vectors  $h$  of test data, if any.

We assume that the analysis has been performed by regression and eigenanalysis, that is, we are given sets of scores  $\Theta$  and sets of fits

$$(22) \quad \eta = B_{\text{OS}}^T h$$

for a given predictor vector  $h$ . It is convenient to translate these directly into the required projections. For example,  $\eta$  might have been produced by a complicated nonparametric regression model (for which we have a software algorithm) with many hidden basis functions; we prefer to avoid these and the innards of the software by simply asking for the fits, or predictions of  $\eta$  at new points.

Starting with the class means, we rewrite (11) using the definition of  $M$  and apply (19),

$$(23) \quad \Theta D_{\alpha} = MB = MB_{\text{LDA}} D_{(1-\alpha^2)^{1/2}},$$

and, denoting with  $\theta^j$  the scores of class  $j$ :  $\Theta = (\theta^1, \dots, \theta^J)^T$ , we get

$$(24) \quad B_{\text{LDA}}^T m^j = \alpha_j / (1 - \alpha_j^2)^{1/2} \theta^j.$$

For actual plotting, one peels off the first two or three dimensions, that is, one uses the first two or three scores for each class; see Figures 4 and 9 for an example.

To map sets of fits  $\eta$  to projected predictors  $B_{\text{LDA}}^T h$ , we only need to apply (21):

$$(25) \quad B_{\text{LDA}}^T h = D_{[\alpha^2(1-\alpha^2)]^{-1/2}} B_{\text{OS}}^T h = D_{[\alpha^2(1-\alpha^2)]^{-1/2}} \eta.$$

Again, for plotting, one picks the first two or three dimensions, that is, the first two or three elements from the set of fits.

In Figure 6 we plotted the penalized class discriminant functions  $\Sigma_W^{-1} m^j$ . From (19) and (36) in the Appendix, we see that  $\Sigma_W^{-1} m^j = B_{\text{LDA}} B_{\text{LDA}}^T m^j$ . Here we need  $B_{\text{LDA}}$  or, equivalently,  $B_{\text{OS}}$ .

*3.6. Distance calculation and classification.* Like many other statistical procedures, LDA can be derived from suitable normal assumptions. In the present paper, these assumptions are not made, but the derivations from them are used as heuristic guides. The usual assumptions for LDA are that the predictor vectors follow multivariate normal distributions with different mean vectors but common covariance matrices among the classes. Assuming (unrealistically) that the parameters in this model are known, classification according to maximum posterior class probability results in a *nearest class mean* rule, where “nearest” is measured in terms of the shared within-class Mahalanobis distance. In practice, the rule is applied with parameters estimated from training data ( $m^j$  and the penalized  $\Sigma_W$  being our choices) and adjustments for unequal class sizes: classify a test predictor  $h$  as class  $j$  if

$$(26) \quad (h - m^j)^T \Sigma_W^{-1} (h - m^j) - 2 \log p_j$$

is minimized over  $j = 1, \dots, J$ . The last term is the class size adjustment ( $p_j = N_j/N$ ), while the crucial first term is the (penalized) squared Mahalanobis distance  $d(h, m^j)$ . Sometimes the sample priors  $p_j$  do not reflect the population priors  $\pi_j$  (e.g., a stratified sample); in this case some external estimate of the  $\pi_j$  should be used.

In order to appreciate the effect of penalization on the metric in  $d(h, m^j)$ , write  $\Sigma_W^{-1}$  as  $(W + \Omega)^{-1}$ , where  $W$  is the unpenalized within covariance. Consider the deviations of two different observations from a centroid:  $O_1$  being close in a “smooth” coordinate, far in a “rough”, and  $O_2$  being the opposite. We would prefer to say that  $O_1$  is closer. However,  $W^{-1}$  will have large eigenvalues for rough directions and will favor  $O_2$ ; the penalty  $\Omega$  will typically have the reverse eigenstructure to  $W$  and so cancels out this effect.

In the Appendix, we first reconstruct the fact that classification can be done based on Euclidean distance in the full space of discriminant variates  $B_{\text{LDA}}^T h$  of  $h$ . Second, we show that the results of optimal scoring, the scores  $\Theta$  and the fits  $B_{\text{OS}}^T h$ , are all that is needed for classification, a fact which we learned from Breiman and Ihaka (1984). This is very useful when we want the regression procedure to select the amount of smoothing, for it shows that the residual sum-of-squares in the regression problem is intimately related to the classification distance [Hastie, Tibshirani and Buja (1994)]. Our proofs are shorter and cover the case of penalized within-class covari-

ances. Third, one can perform classification using distances in the  $K < J - 1$  reduced-dimensional discriminant subspace by using the corresponding reduced set of regression fits and scores.

These results establish that the original class-adjusted Mahalanobis distance (26) is equivalent to each of the following distance measures (when used in a relative sense for classification):

1.  $d(h, m^j) - 2 \log p_j$ ;
2.  $\|B_{\text{LDA}}^T(h - m^j)\|^2 - 2 \log p_j$ ;
3.  $\|D_{\alpha(1-\alpha^2)^{1/2}}(\eta - \bar{\eta}^j)\|^2 - 2 \log p_j$ ;
4.  $\|D_{(1-\alpha^2)^{-1/2}}(\theta^j - \eta)\|^2 - \|\theta^j\|^2 - 2 \log p_j$ ;
5.  $\|D_{(1-\alpha^2)^{-1/2}}(\theta^j - \eta)\|^2 - 1/p_j - 2 \log p_j$ .

Among these, only the last cannot be used in a dimension-reduction mode since it relies on the presence of  $J - 1$  discriminant coordinates. In the third expression,  $\bar{\eta}^j = B^T m^j = D_\alpha^2 \theta^j$ .

**4. Model considerations.** Both the examples in this paper arise from discretized analog signals: cases where the data can be viewed as functions, sampled for computer representation. In the first example, the spectrum for an utterance is a function of frequency; in the second, the picture of the handwritten digit is a function of two spatial coordinates. As such the data for the  $i$ th sample can be interpreted as a discretized realization  $h_{i,k} = h_i(s_k)$  of a stochastic predictor process  $h(s)$  evaluated at discrete values  $s_k$  of a (continuous) index domain. We assume (a) that the  $j$ th class is observed with relative frequency  $\pi_j = P[G = j]$  and that (b) its predictor process  $h(s)$  has a class-specific mean function  $\mu_j(s) = E[h(s)|G = j]$  and (c) a covariance function  $\Sigma(s, t) = E\{[h(s) - \mu_j(s)][h(t) - \mu_j(t)]|G = j\}$  that is shared among the classes.

This provides us with a model underlying the methods described in this paper. It is natural to think of a functional version of LDA [Kiiveri (1992), Section 1.2] in terms of this model. Ramsay and Dalzell (1991) coined the name *functional data analysis* for this kind of problem.

The functional canonical variate problem in discriminant analysis consists of finding normalized functions  $\beta_k(s)$  such that the associated functionals  $\eta_k = \int \beta_k(s)h(s) ds$  have means that are optimally separated among the classes, that is,  $\sum_{j=1}^J \pi_j E(\eta_k|G = j)^2$  is maximized. The normalization has the form  $\iint \beta_k(s)\Sigma(s, t)\beta_l(t) ds dt = \delta_{kl}$ .

For the functional classification problem in discriminant analysis, one assumes that the predictor processes are Gaussian. The Bayes optimal classifier assigns a predictor function  $h(s)$  to the class  $j$  that minimizes  $d(h, \mu_j) - 2 \log \pi_j$ , where the Mahalanobis distance  $d$  is defined as

$$d(h, \mu_j) = \iint [h(s) - \mu_j(s)] \Sigma^{-1}(s, t) [h(t) - \mu_j(t)] ds dt.$$

For this to be meaningful, one has to assume that an inverse  $\Sigma^{-1}$  of the covariance operator  $\beta(\cdot) \mapsto \int \Sigma(\cdot, t)\beta(t) dt$  exists, but even if it exists, it generally cannot be represented by a kernel  $\Sigma^{-1}(s, t)$ . The double integral in the definition of  $D(h, \mu_j)$  therefore has to be taken with a grain of salt.

For estimation in the functional model, we have to consider two levels of asymptotics:

1. we obtain  $N$  realizations  $h_i$  of the predictor process (sample size  $N$ );
2. we observe the realizations at  $m$  discrete locations  $s_k$  (resolution  $m$ ).

Leurgans, Moyeed and Silverman (1993), in the context of canonical correlation analysis, focus on the first case and essentially view the second as a computational approximation of integrals by sums. For infinite-dimensional predictors  $h(s)$  ( $m = \infty$ ), they show that in order to obtain consistent estimates of  $\beta(s)$  it is essential to regularize the sample estimate of the covariance function (see their Propositions 1 and 2). This is intuitively obvious, for, with infinite resolution  $m$  and finite sample size  $N$ , we have always more variables than observations, and well-known degeneracies occur. They then derive the rates at which the amount of regularization should decrease to zero as  $N$  grows large in order to achieve consistent estimates of canonical variates. These results carry over to discriminant analysis with minor modifications: in their functional canonical variates problem, specialize the second process to a stochastic indicator vector indicating class membership, and remove penalization for this  $J$ -dimensional “class membership process.”

Estimation of the Bayes optimal classification criterion poses an equally obvious problem: any estimate  $\hat{\Sigma}(s, t)$  of the covariance function  $\Sigma(s, t)$  will have rank no greater than the sample size  $N$ . It is therefore impossible to invert the estimated covariance operator in order to calculate the Mahalanobis distance, unless it is regularized with a suitable penalty. There is no necessity, however, to develop asymptotic theory for this form of the classification criterion since other equivalent forms of the criterion are based on canonical variates (Section 3.6) and hence require inversion of the covariance operator only on a finite-dimensional subspace.

As a final remark, one might argue that asymptotic theory of estimation in a functional framework is not realistic since for all practical purposes a finite resolution has to be chosen. The argument could continue along the lines that the resolution should be selected commensurate with the sample size, and a realistic asymptotic theory should determine at what rate the resolution  $m$  can be increased as a function of the sample size  $N$  in order to achieve consistency. Such an approach would essentially use resolution as a regularization parameter. The problem with this approach is that the choice of resolution is and should be guided by computational feasibility rather than sample size: if we face a very small  $N$ , it might be feasible to choose a rather large  $m$  without running into problems of computer time and memory exhaustion. Low resolution amounts to throwing away data—a rather crude regularization method. It is more informative to make use of as much data as possible and derive regularization from substantive arguments, such as smoothness considerations on the index domain.

**5. Example: smooth canonical functions for classifying log-periodograms.** The following analysis is taken from joint work of Andreas Buja, Werner Stuetzle and Martin Maechler. The data in this example were extracted from the TIMIT database (TIMIT Acoustic-Phonetic Continuous Speech Corpus, NTIS, U.S. Department of Commerce), which is a widely used resource for research in speech recognition. We formed a small test problem by selecting five phonemes for classification based on digitized speech from this database. The phonemes are transcribed as follows: “sh” as in “she”; “dcl” as in “dark”; “iy” as the vowel in “she”; “aa” as the vowel in “dark”; and “ao” as the first vowel in “water.” From continuous speech of 50 male speakers, we selected 1000 speech frames of 32 ms duration, approximately 2 examples of each phoneme from each speaker. Each speech frame is represented by 512 samples at a 16 kHz sampling rate, and each frame represents one of the above five phonemes. The breakdown of the 1000 speech frames into phoneme frequencies is as follows:

sh	dcl	iy	aa	ao
191	167	269	147	467

From each speech frame, we computed a log-periodogram, which is one of several widely used methods for casting speech data in a form suitable for speech recognition. Thus the data used in what follows consist of 1000 log-periodograms of length 256, with known class (phoneme) memberships.

Figure 3 shows a sample of 10 log-periodograms in each phoneme class. It is known that periodograms have rather erratic statistical properties; the figure supports this. In order to turn them into reasonable estimates of an underlying spectral density function, a certain amount of smoothing is required. However, the interest in speech recognition is not faithful estimation of spectral densities but discrimination among speech units such as phonemes. We therefore envision a role for smoothing at the level of linear *functionals* on log-periodograms rather than smoothing the periodograms themselves.

If we think of log-periodogram vectors as functions  $h(f)$  evaluated at finitely many frequencies  $f$ , we have the task of estimating linear functionals  $\eta(h) = \int_f \beta(f)h(f) df$ , approximated by  $\sum_f h(f_j)\beta_j$ , that best discriminate between the phonemes. The representers  $\beta(f)$  and their discretized versions (the sets of linear coefficients) should be estimated in such a way that they depend smoothly on the frequency  $f$ , for reasons presented in Sections 1 and 4. This can be achieved in an obvious way, for example, with a second derivative penalty on the coefficients  $\beta(f)$ :  $\lambda \int [\beta''(f)]^2 df$  [Wahba (1990)]. The penalty that we actually used for the present speech problem has the form  $\lambda \int [\beta''(f)]^2 w(f) df$ , where  $w(f)$  was chosen to penalize the higher frequencies more. This is suggested by the experimental observation that the capability of the human ear to discern acoustic detail decreases rapidly for the frequencies above 1 kHz. In speech recognition, it is standard to use techniques with higher resolution in lower frequencies. Our weight function  $w(f)$  is constant for values of  $f$  up to 1 kHz, and decreases linearly there-

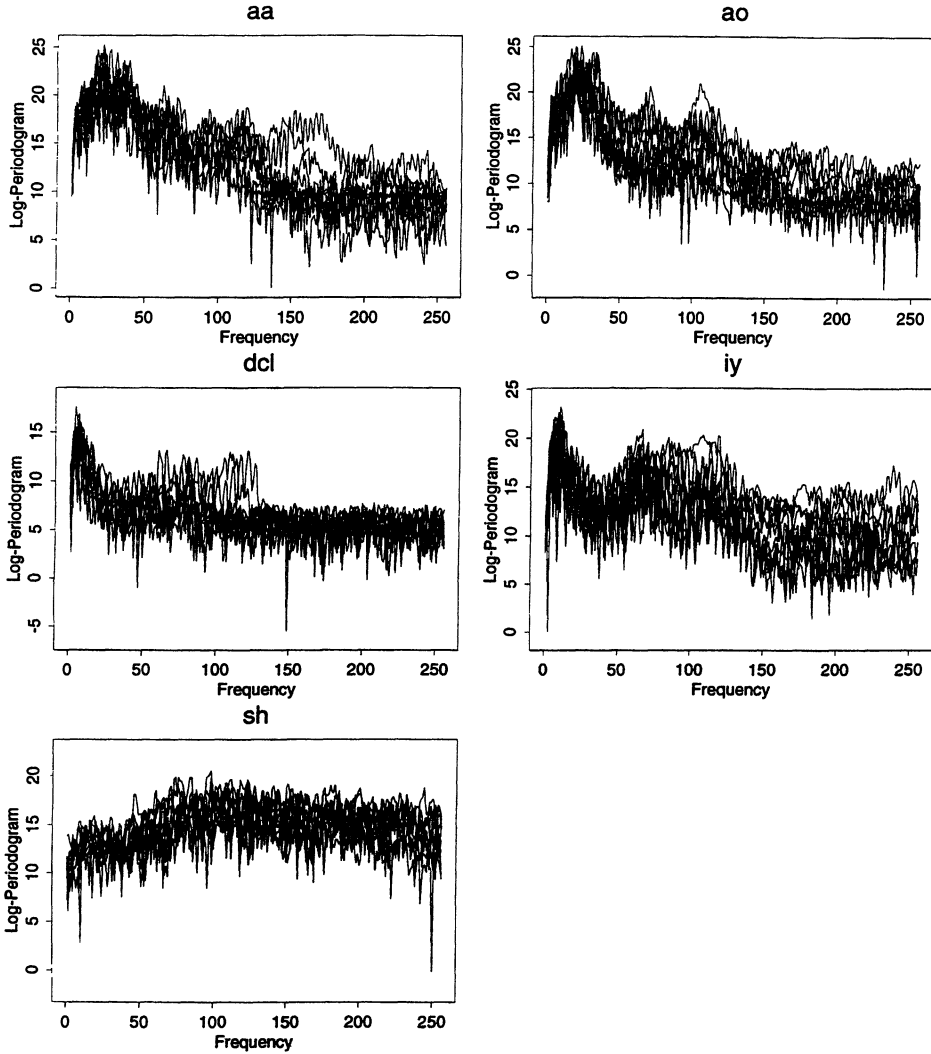


FIG. 3. Each plot shows a sample of 10 log-periodograms within each phoneme class. The periodograms consist of a smooth trend plus high-frequency oscillations.

after. We refer to this as an *improper spline* penalty; further details are given in Section 7.

In Figure 1 in Section 1 we show the four canonical LDA coefficients  $\beta(f)$  plotted as functions of a frequency scale (the plotting range 0–256 maps to 0–8 kHz). The jagged curves represent the raw LDA coefficients, while the smooth curves show penalized coefficients.

The overwhelming impression from these plots is of course the extremely wild behavior of the unpenalized LDA coefficients as compared to the penal-



ized ones. Except for the very lowest frequencies, there is barely any structure that could be discerned by eye in the unpenalized curves. Although not visually evident in Figure 3, there are positive correlations among log-periodogram values over large neighborhoods on the frequency scale. In fact, among more than 32,000 (256 choose 2, to be exact) pairwise within-class correlations of log-periodogram values, fewer than 5% were negative. For the reasons given in Section 1, this leads to negatively correlated coefficients.

The penalized curves shown in Figure 1 are easily recognized by their smoothness. In terms of fitted degrees of freedom, the reduction in the penalized coefficients is dramatic: while the unpenalized coefficients have 256  $df$ , the penalized coefficients have only 30  $df$ . (We computed  $df$  as the trace of the smoother matrix in the equivalent optimal scoring problem; further details in Section 7.) The smoothness of the penalized curves allows interpretation of the coefficients as contrasts that set various frequency ranges against each other. The action is mostly in the low frequencies, as expected. Figure 4 shows that the first penalized discriminant coordinate function sets the phonemes *aa* and *ao* apart from the phonemes *dcl*, *iy* and *sh*. According to Figures 1 and 3, this coordinate relies on a peak (so-called formant) between 500 and 1000 Hz (roughly frequencies 16 to 32) in the log-periodograms for *aa* and *ao*. According to Figure 4, the second penalized discriminant coordinate sets the phoneme *sh* against the phonemes *dcl* and *iy*. A look at Figures 1 and 3 shows that this contrast picks up a peak near zero frequency in the log-periodograms for *dcl* and *iy* which is absent for *sh*.

So far, we discussed only qualitative features of penalized discriminant analysis for this data example, but the bottom-line numbers, namely, misclassification rates, are favorable for penalization as well.

In order to assess test sample performance, we conducted a simulation study using an additional 387 speakers. From the combined set we randomly selected 50 speakers, fitted both the LDA model and a variety of PDA models at different levels of regularization and used them to classify the remaining examples. This was repeated 50 times, and the results are shown in Figure 5. In the lower left panel we see boxplots of the test error for the improper spline as a function of  $df$ , with the rightmost entry corresponding to the unpenalized LDA. The minimum median of 0.073 is achieved around 30  $df$ , and all are lower than that for unpenalized LDA (median 0.086). The amount of regularization needed depends on the size of the training set, as well as the purpose of the model (we might favor smoother coefficient functions for interpretability). The figure suggests that the test error rates of PDA are relatively insensitive to the amount of smoothing: any reasonable amount between 20 and 80  $df$  does considerably better than no smoothing. By comparison, the top left panel shows the corresponding error rates on the training data. As expected, LDA overfits and does the best, while PDA achieves a median value slightly lower than the 0.073 it achieved on the test data.

Included in Figure 5 are the results of some other approaches to PDA. The second column of figures was based on a hand-crafted filtering approach

## Canonical Variate Plot --- Phoneme Training Data

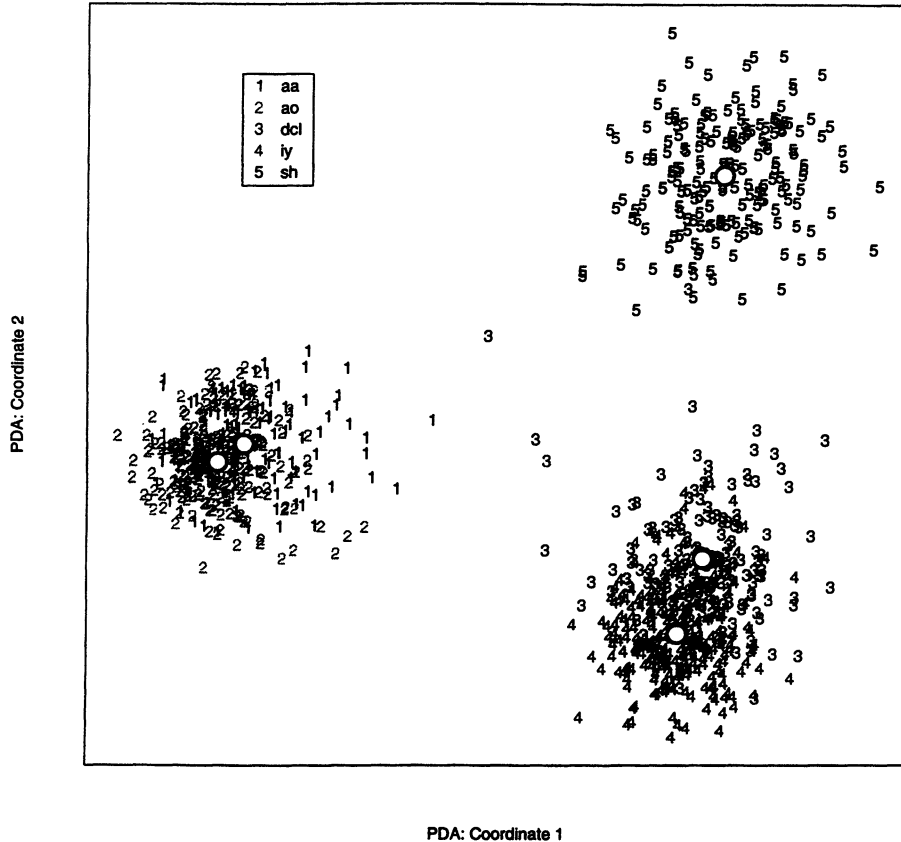


FIG. 4. *The first two penalized canonical variates, evaluated for the phoneme training data. The true class identities are indicated. As might be expected, the two vowel sounds “ao” and “aa” are confused. The circles indicate the class centroids.*

popular in speech research, which uses piecewise linear basis functions. The knot spacings are chosen to be uniform up to 1 kHz, and then increasing logarithmically. Our piecewise linear weighting function  $w(f)$  used with the improper spline penalty was chosen to correspond empirically with this choice of knots. These filter bases are used to create linear combinations of the original predictors and are then used in place of them. The misclassification results are comparable to those for the improper spline.

The third column was obtained by using an unstructured ridge penalty to regularize, which amounts to shrinking the within covariance matrix toward a scalar covariance matrix. Again the misclassification results are only slightly worse than for the first two columns. All three approaches can be viewed as versions of PDA, enforcing spatial smoothness of the coefficient

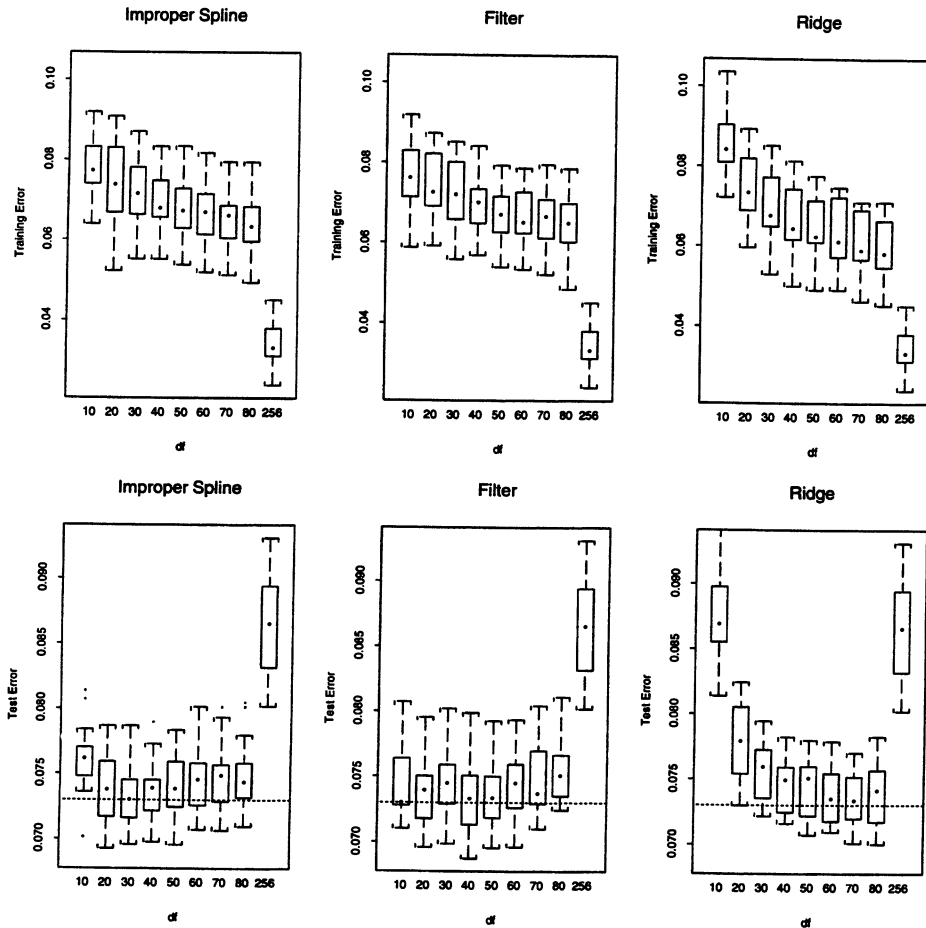


FIG. 5. Training and test errors as a function of  $df$  for three versions of PDA, compared with the unrestricted LDA ( $df = 256$ ). The three versions are (left) regularized using a weighted smoothing penalty, (middle) asymmetric filters and (right) simple ridging. All three show comparable classification performance at their optimum  $df$ .

functions. The filter approach explicitly uses piecewise smooth bases to represent the coefficients. Although ridging appears to treat all coefficients equally, one can show that those contrasts corresponding to small eigenvalues of the total covariance are shrunk more than those corresponding to the large eigenvalues. Since the positive autocovariances lead to longer-trend dominant eigenfunctions, to some extent the ridging shrinks toward similarly behaved coefficient functions [Hastie and Mallows (1993)]. The improper-spline approach, on the other hand, explicitly biases toward smooth coefficient functions.

Figure 6 shows the *class discriminant functions* for the five classes, comparing LDA, our improper-spline PDA and ridge PDA, each at their

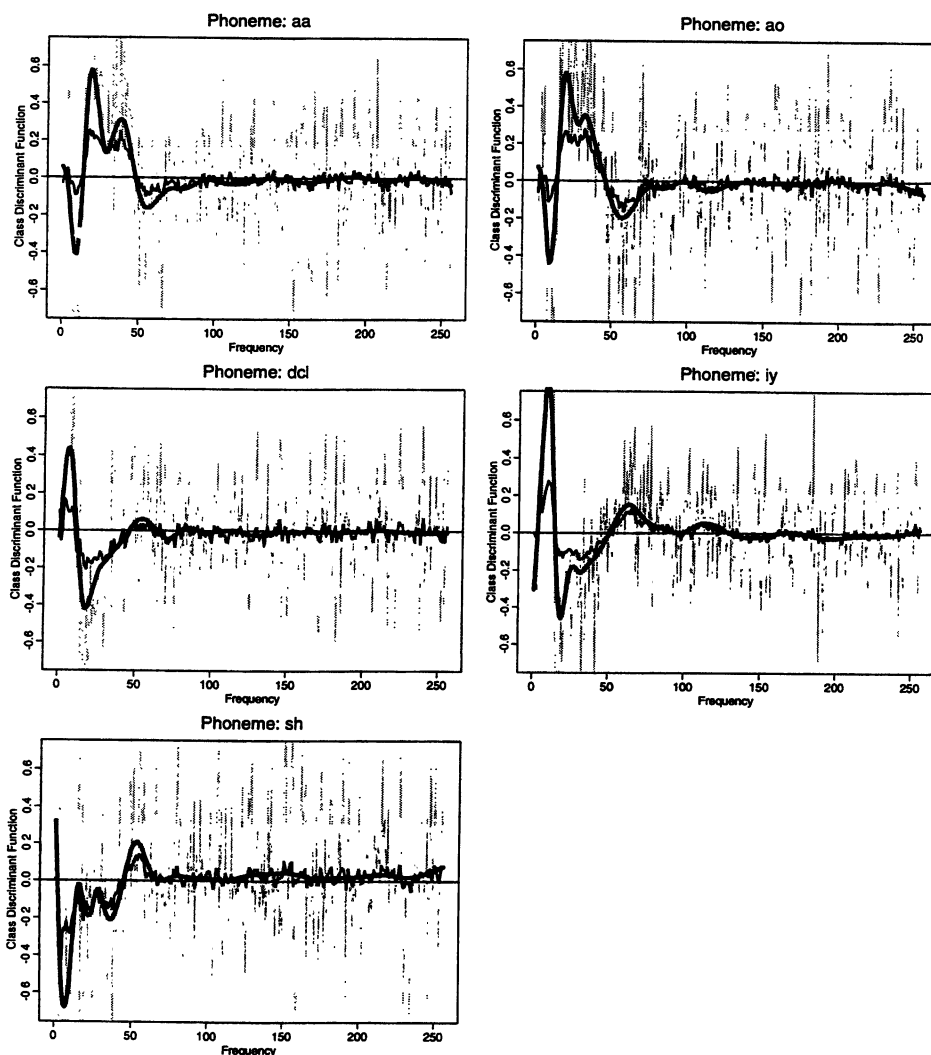


FIG. 6. Class discriminant functions for the five phonemes: these are analogs of the functions  $\Sigma_W^{-1} m^j$  for linear discriminant analysis, except the within covariance here is penalized. The thick solid curve represents PDA using the improper spline penalty and 30  $df$ . The raw LDA coefficients are in the background. The wiggly curve is the ridge version of PDA, using 70  $df$ .

optimal  $df$  according to Figure 5. These are the equivalent versions of  $\Sigma_W^{-1} m^j$  for linear discriminant analysis (see the last paragraph in Section 3.5). These are distinguishing contrast coefficients for each phoneme. As expected the two “a” vowels are very similar. The ridge curve is very wiggly, but still manages to portray the main features. Since ridge shrinks toward the origin, far less shrinking was asked for (70  $df$ ) in order to allow the main features to appear; the improper spline, on the other hand, shrinks toward smoothness and could

afford smoother functions. For space reasons we do not display the curves for the filter method, which look similar to those for the improper spline, albeit piecewise linear.

**6. Example: image analysis and character recognition.** A current “hot” topic in the field of pattern recognition is the automatic reading of handwritten addresses and zip-codes from envelopes. Le Cun, Boser, Denker, Henderson, Howard, Hubbard and Jackel (1990) implemented a successful procedure which included as a primitive the classification of isolated handwritten digits. We used a subset of their data for training and testing, a sample of which is shown in Figure 7.

Le Cun et al. normalized the binary images for size and orientation, resulting in 8-bit,  $16 \times 16$  gray-scale images. They used the 256 pixel values as the inputs to a multilayer-neural-network model, which is trained to classify the images. They report overall misclassification rates on test data of under 5%. The goal in this section is not to compete with their highly tuned procedure, but to show that a traditional statistical method, LDA, can be improved by spatial regularization. One could expect that any classification procedure that relies on coefficients at a pixel level (including neural networks) would similarly benefit from regularization.

Our approach here was to use the same normalized data as Le Cun and coworkers (kindly supplied by J. Bromley), but to use simpler and more classical approaches to discrimination. We used the first 2000 images as training data, and the following 2000 as a preliminary validation set.

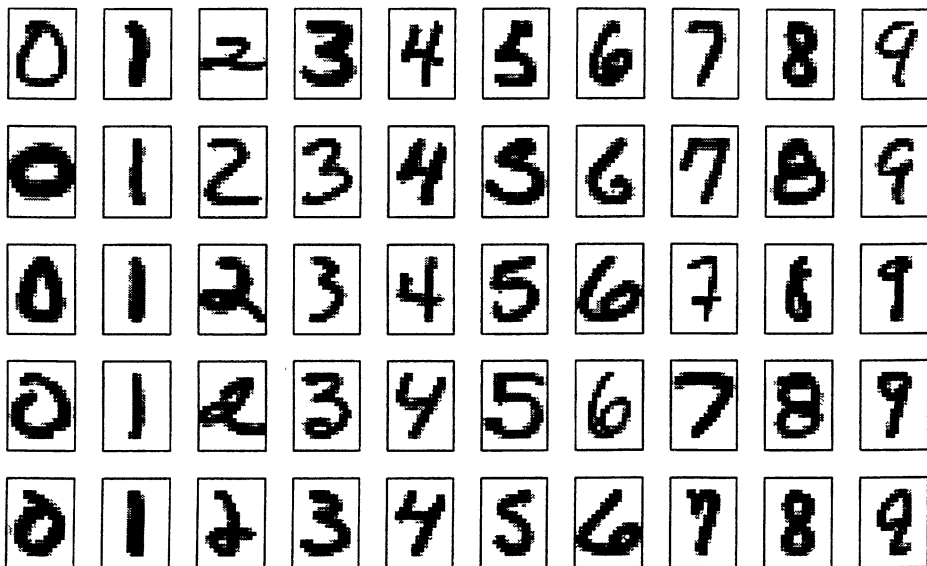


FIG. 7. A random selection of digitized handwritten digits: each image is an 8-bit, gray-scale version of the original binary image, size and orientation normalized to  $16 \times 16$  pixels.

As a first step we fitted a standard LDA model to the training data. The misclassification rate was 3.1% on the training data, but 11% on the validation set. The training error is less than a third of the test error and suggests we may be in an overfit situation, despite the simplicity of the technique. Indeed, the nine discriminant functions each have 256 coefficients, and despite the normalization constraints, resulted in over 2000 independent parameters!

Apart from the initial normalization, our procedure did not take advantage of the spatial correlation in the data. Even though we have 256 pixels per image—ostensibly 256 independent pieces of information—this spatial correlation suggests the real number is far less.

Figure 2 in Section 1 consists of nine pairs of images; the left member of each pair represents the LDA coefficients as images, with the hope of gaining some insight into the important contrasts found. These salt-and-pepper images reveal no structure and once again reflect the correlation between neighboring pixels, resulting in a strong negative correlation between coefficients.

We fit a PDA model using a Laplacian penalty  $J(\beta)$  to constrain the coefficients to be spatially smooth. Denote the coefficients as functions  $\beta(x, y)$ , where  $x$  and  $y$  refer to horizontal and vertical coordinates in the two-dimensional image. Then

$$(27) \quad J(\beta) = \int_{x,y} \left[ \frac{\partial^2 \beta}{\partial x^2} + \frac{\partial^2 \beta}{\partial y^2} \right]^2 dx dy,$$

where the term within square braces is the Laplacian. Further details are given in Section 7.

The right member of each pair in Figure 2 represents the PDA fit, where we chose the smoothing parameter to achieve 40  $df$  for each coordinate. The misclassification rates are now 6.1% for the training data, and 8.2% for the test data; a 25% reduction in validation error using 85% fewer parameters!

In order to validate these results, we conducted a simulation study similar to that used in the previous example. We randomly sampled 2000 images from the combined set of 4000 as training images and used the remainder as test images. This was repeated 50 times, and the results for several methods are summarized in Figure 8. We have used about 40  $df$  for all the examples. We tried different amounts of smoothing and found that for less  $df$ , the test results uniformly got worse quite rapidly, and for more  $df$  they remained the same or got worse very slowly. From these results it seems our initial division into training and test set yielded slightly pessimistic results.

The smoothness not only gives us improved misclassification rates, but offers interpretability as well. The first PDA coefficient in Figure 2 looks like a white 0 with a vertical center darker than the remainder of the image. This contrast image can be expected to yield positive scores for 1's and negative scores for 0's. Figure 9 confirms this, and as we might expect 7's and 9's are on the same side as 1's, while 6's look more like 0's. The second PDA

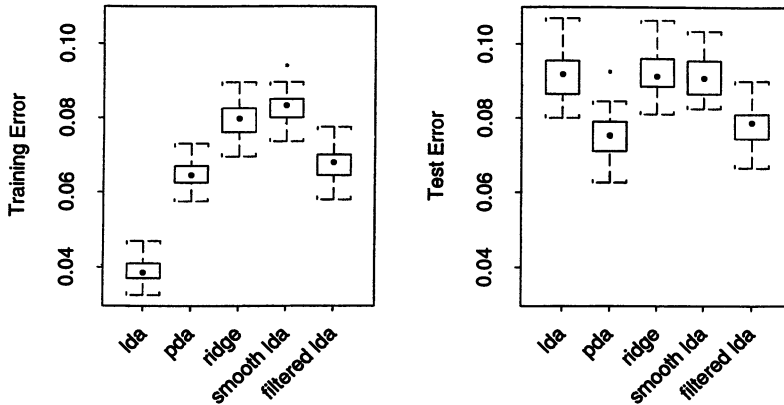


FIG. 8. Each boxplot shows the misclassification results for 50 resampling trials, the left panel showing training or resubstitution error, the right panel showing test error. All methods use about 40 degrees of freedom. The method labelled “smooth LDA” was obtained by smoothing the LDA coefficients as a function of spatial coordinates, using the same Laplacian penalty as the PDA fits. The method labelled “filtered LDA” corresponds to LDA on a 44-dimensional set of filter coefficients, described in the text.

coefficient in Figure 2 gives a positive score to images dark in the bottom left corner, and a negative score to images dark in the bottom middle and top left; 6’s fall into the first category and 7’s and 9’s into the second, and both are confirmed in Figure 9. The interpretation of the coefficients becomes more difficult for the lower-order coefficients.

For the filtering approach we used a tensor-product basis of polynomials in each of the spatial coordinates, with total degree restricted to 8 and thus  $m = 44$  basis functions. The images derived for this and the other methods were not visually informative so we have omitted them.

**7. Penalized regression methods.** In this section we give more details on the different forms of penalized regression outlined in Section 1, used in the examples and appearing in equations (5) and (4). The response is a scored version of  $G$ , thus a scalar, numeric variable, which for simplicity we refer to as  $Y$ , with realizations  $y_i$ .

Ridge regression is the oldest and simplest method for regularizing linear regression problems with nearly collinear predictors. It is trivially a penalized least squares method since the ridge estimator  $\hat{\beta} = (H^T H + \lambda I)^{-1} H^T y$  minimizes the criterion

$$\|y - H\beta\|_N^2 + \lambda \|\beta\|_m^2,$$

where  $N$  is the sample size and  $m$  the number of predictor variables. The matrix  $\lambda I$  and its associated penalization  $\lambda \|\beta\|_m^2$  are adequate when neither prior knowledge nor purpose of investigation dictate a more adapted kind of

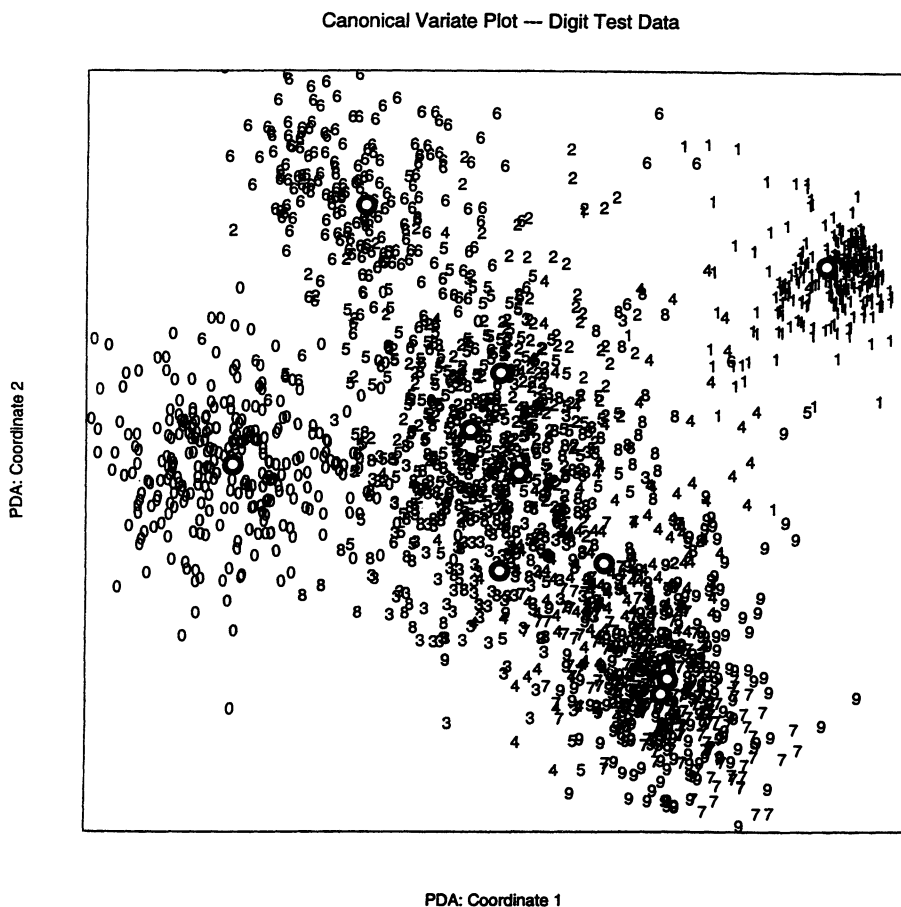


FIG. 9. *The first two penalized canonical variates, evaluated for the test data: the circles indicate the class centroids. The first coordinate contrasts mainly 0's and 1's, while the second contrasts 6's and 7/9's.*

regularization. More generally the estimator has the form  $\hat{\beta} = (H^T H + \lambda \Omega)^{-1} H^T y$ , and it minimizes

$$(28) \quad \|y - H\beta\|_N^2 + \lambda \beta^T \Omega \beta.$$

Here  $\Omega$  is a more structured penalty matrix and imposes smoothness with regard to an underlying space, time or frequency domain.

We also want to consider smoothness at the level of functions  $\beta(s)$  on this domain, as in Section 4. The functional discriminant analysis model translates via optimal scoring into the functional regression model  $\eta = \int \beta(s) h(s) ds$  for the mean of  $Y$ , and an appropriate criterion might be

$$(29) \quad E_Y (Y - \eta)^2 + \lambda J(\beta)$$



for some penalty functional  $J$ , a seminorm of the space of functions being considered. In discretizing the problem, we have a sample of responses  $y_i$  and functions  $h_i$  measured at a set of  $m$  values of  $s$ , but we can leave  $\beta(s)$  as a function, chosen to minimize

$$(30) \quad \sum_{i=1}^N \left( y_i - \sum_{l=1}^m h_l(s_i) \beta(s_l) \right)^2 + \lambda J(\beta).$$

We now focus on three particular applications: the two examples given in this paper, and the nonparametric regression approach used in FDA.

7.1. *Speech data: improper splines.* Here the penalty has the form

$$(31) \quad J(\beta) = \int [\beta''(f)]^2 w(f) df,$$

where  $w(f)$  allows the penalty to give more weight to the lower frequencies. If  $w$  were missing, the solution to (31) would be a natural cubic-spline [Wahba (1990)] with knots at the sampled frequencies  $f_i$ . With  $w$  present, the solution will not in general be a spline, and the character of the solution, if it exists, will depend on  $w$ . We sidetrack these issues by adapting the spline penalty to approximate (31), and refer to our solution as an *improper spline*. With  $w$  absent, the  $m$ -dimensional spline solution can be parameterized by the values  $\beta_l = \beta(f_l)$  themselves, and the penalty can be written as  $J(\beta) = \beta^T \Omega \beta$ , with  $\Omega = \Delta^T C^{-1} \Delta$  [Green and Yandell (1985), for example]. Here  $\Delta$  is a second-difference operator for approximating second derivatives, and  $C^{-1}$  can be viewed as a kernel for approximating the integral. If  $D_w$  is a diagonal matrix with entries  $w(f_l)$ , we use  $\Omega_w = \Delta^T D_w^{1/2} C^{-1} D_w^{1/2} \Delta$  as the penalty in (28).

7.2. *Image data: tensor product splines.* Here our coefficients  $\beta(x, y)$  are indexed by  $x$  and  $y$ : horizontal and vertical coordinates in the two-dimensional image, and  $h(x, y)$  are gray-scale values. We used the Laplacian penalty

$$(32) \quad J(\beta) = \int_{x,y} \left[ \frac{\partial^2 \beta}{\partial x^2} + \frac{\partial^2 \beta}{\partial y^2} \right]^2 dx dy.$$

Since  $x$  and  $y$  each assume 16 uniformly spaced values on a lattice in  $R^2$ , for computational simplicity we use the discretized version of  $J(\beta)$  described in O'Sullivan [(1991), page 635]. This uses the discrete approximation  $J(\beta) = \beta^T \Omega \beta$ , where the  $256 \times 256$  penalty matrix  $\Omega = \Delta^T \Delta$  is based on the discrete approximation to the Laplacian:  $\Delta = D_x \otimes I_y + I_x \otimes D_y$ . The 256-vector  $\beta$  refers to the vectorized version of the  $16 \times 16$  matrix of coefficients (in row order). Here  $D_x$  and  $D_y$  are each  $16 \times 16$  matrices that approximate a second derivative by second differences. O'Sullivan (1991) gives an explicit Fourier representation of the spectral decomposition  $\Omega = \Gamma D_\beta \Gamma^T$ . This simplifies the solution of (28), since we can transform to new coordinates for which the penalty is diagonal.

**7.3. Nonparametric regression.** The regularized regressions we have seen so far penalize the coefficients of the predictors for roughness (in some spatial domain). In the form of nonparametric regression considered here, we first expand the predictors into a large set of basis functions, and then restrict their coefficients to be smooth as a function of the predictors themselves.

For the purposes of the discussion here, we focus on nonparametric additive models for the regression function:  $\eta(x) = \sum_{k=1}^p f_k(x_k)$  [Buja, Hastie and Tibshirani (1989), Hastie, Tibshirani and Buja (1994)]; and an appropriate penalized least-squares criterion

$$(33) \quad \sum_{i=1}^N (y_i - \eta(x_i))^2 + J(\eta; \lambda),$$

with

$$J(\eta; \lambda) = \sum_{k=1}^p \lambda_k [f_k''(u)]^2 du.$$

The solution can be represented in terms of a finite-dimensional basis  $h(x)$ :  $\eta(x) = h(x)^T \beta$ . Here  $h(x)$  is a union of spline basis functions for each coordinate function. Suppose  $h$  has  $m$  components, and let  $H$  be the  $N \times m$  basis matrix with  $i$ th row  $h(x_i)$ ;  $J(\eta; \lambda)$  can be expressed as a quadratic form in  $\beta$ :  $J(\eta; \lambda) = \beta^T \Omega \beta$ , where  $\Omega$  is a  $m \times m$  block-diagonal, nonnegative definite penalty matrix. This again has the form (28).

The number  $m$  of basis function can be quite large; for a univariate smoothing spline there are as many basis functions as there are unique values of  $x_i$  (potentially  $N$ ). The additive-spline model can have up to  $Np$  basis functions! Fortunately there are efficient iterative algorithms for computing the solutions, especially for additive spline models. More details on the additive spline model (and several examples using them in this context) are given in Hastie, Tibshirani and Buja (1994).

**7.4. Degrees of Freedom.** In the examples we have referred to the effective degrees of freedom of a smooth fit. This is a useful translation of the smoothing parameter into a more meaningful parameter that can be used to calibrate a variety of different linear methods. The vector of fitted values for  $\eta$  in (28) is given by

$$(34) \quad S(\lambda)y = H(H^T H + \lambda\Omega)^{-1} H^T y,$$

where  $S(\lambda)$  is an  $N \times N$  linear operator matrix. The degrees of freedom are defined to be  $df(\lambda) = \text{tr} S(\lambda)$  [Buja, Hastie and Tibshirani (1989), Wahba (1990)]. The function  $df$  is monotone in  $\lambda$ , and in practice we can fix  $df$  and determine the appropriate value of  $\lambda$  with minimal additional computational cost. In the examples in this paper we avoided automatic selection of regularization parameters. Instead we have favored an empirical approach of examining discriminant functions and misclassification performance for a few values of  $df$  and selecting the most informative values.

**8. Summary.** Penalized discriminant analysis appears to be a useful tool for problems such as speech and image recognition, featuring a large number of correlated inputs. An attractive feature of the method is the potential to extract and interpret a small number of discriminant functions and the resulting benefits for scientific understanding of the data. Potential medical applications include the disease classification of pap smears and mammograms.

## APPENDIX

### Details of the distance calculation.

*A.1. Reduction of distances to discriminant coordinates.* The decomposition in (8) is in terms of at most  $J - 1$  eigenvectors of  $\Sigma_{\text{Bet}}$ ; the  $J$ th is trivial with eigenvalue 0 if  $H$  is centered, or else equally trivial with eigenvalue 1 if  $H$  is not centered but includes the constant column. Assuming that  $J < m$ , it is useful to pad out the decomposition to the full  $m$  dimensions: (a) let  $B^* = (B : B_\perp)$  be of size  $m \times m$ , such that  $B^{*T} \Sigma_{22} B^* = I_m$ , and (b)  $D_\lambda = D_{(\alpha^2; 0)}$ , where  $B$  is now the full  $(J - 1)$ -dimensional nontrivial solution and  $D_{\alpha^2}$  is the corresponding nonzero eigenvalues. The following facts can be easily verified:

$$(35) \quad \begin{aligned} \Sigma_{22}^{-1} &= B^* B^{*T} \\ &= BB^T + B_\perp B_\perp^T; \end{aligned}$$

$$(36) \quad \begin{aligned} \Sigma_W^{-1} &= B^* (I - D_\lambda)^{-1} B^{*T} \\ &= B(I - D_\alpha^2)^{-1} B^T + B_\perp B_\perp^T; \end{aligned}$$

$$(37) \quad \Sigma_{\text{Bet}} = B^{*-T} D_\lambda B^{*-1}.$$

An implication is that, confined to the subspace spanned by  $M$ , the (penalized) Mahalanobis distances differ using the metrics  $\Sigma_{22}^{-1}$  and  $\Sigma_W^{-1}$ , but are the same in the orthogonal subspace.

The following equation shows that classification can be based entirely on discriminant coordinates:

$$(38) \quad d(h, m^j) = \|B_{\text{LDA}}^T (h - m^j)\|^2 + \|B_\perp^T h\|^2.$$

The first term follows from (36) together with the relation between  $B_{\text{LDA}}$  and  $B$  from (19). From (37) we see that  $B_\perp^T \Sigma_{\text{Bet}} = 0$ , but since  $\Sigma_{\text{Bet}} = M^T \Sigma_{11} M$ , it follows that  $B_\perp^T M^T = 0$ , and hence  $B_\perp^T m^j = 0 \forall j$ . Thus (38) follows.

If the dimension  $m$  of the (expanded) predictors  $h$  is much larger than the number of classes,  $J$ , the reduction to discriminant coordinates may give considerable savings.

*A.2. Reduction of distances to optimal scoring.* In Section 3.5, we established links between projections based on LDA and OS. In terms of distance

calculations, they imply that classification can be based on scores  $\theta^j$  and fits  $\eta$  rather than class centers  $m^j$  and predictors  $h$ . This can be convenient in settings where the dimension  $m$  of  $h$  is very large; for example, when  $\eta$  is modeled by an additive spline,  $m \sim Np$ , where  $p$  is the number of original predictors in  $x$ . Algorithms for fitting additive spline models efficiently compute the additive fit  $\beta_{\text{OS}}^T h$ ; it is therefore useful to express the classification criterion in terms of this regression as well. To this end, we reformulate (38) based on (24) and (25):

$$\begin{aligned}
\|B_{\text{LDA}}(h - m^j)\|^2 &= h^T B_{\text{LDA}} B_{\text{LDA}}^T h - 2h^T B_{\text{LDA}} B_{\text{LDA}}^T m^j \\
&\quad + m^{jT} B_{\text{LDA}} B_{\text{LDA}}^T m^j \\
(39) \quad &= \eta^T D_{\alpha^{-2}(1-\alpha^2)^{-1}} \eta - 2\eta^T D_{(1-\alpha^2)^{-1}} \theta^j \\
&\quad + \theta^{jT} D_{\alpha^2(1-\alpha^2)^{-1}} \theta^j \\
&= \eta^T D_{(\alpha^{-2}-1)(1-\alpha^2)^{-1}} \eta + (\theta^j - \eta)^T D_{(1-\alpha^2)^{-1}} (\theta^j - \eta) \\
&\quad - \theta^{jT} D_{(1-\alpha^2)(1-\alpha^2)^{-1}} \theta^j \\
&= \|D_{\alpha^{-1}} \eta\|^2 + \|D_{(1-\alpha^2)^{-1/2}} (\theta^j - \eta)\|^2 - \|\theta^j\|^2.
\end{aligned}$$

This result in itself proves the assertion. A closer look at the derivation reveals that it is of greater generality than meets the eye. Assume one wishes to perform dimension reduction and base the distance calculations for classification on the first  $K < J - 1$  discriminant coordinates only. One would then use the following:

1.  $B_{\text{LDA}}$  reduced to the first  $K$  columns;
2.  $\theta^j$  reduced to the first  $K$  scores;
3.  $\eta$  reduced to the first  $K$  fitted values.

It turns out that the derivations of (39) still hold. Thus, scores and fitted values are a base for distance calculation at any level of dimension reduction.

Assuming that no dimension reduction is desired, one can continue and express (39) in a variety of ways, for example, by rewriting the equation (38) for the Mahalanobis distance with the help of (35):

$$(40) \quad d(h, m^j) = h^T \Sigma_{22}^{-1} h + \|D_{(1-\alpha^2)^{-1/2}} (\theta^j - \eta)\|^2 - \|\theta^j\|^2.$$

A final touch concerns the term  $\|\theta^j\|^2$ : Let  $\Theta'$  denote  $\Theta$  augmented with the trivial column  $\theta_j = 1$ . Then  $\Theta'$  is square and nonsingular with  $\Theta'^T \Sigma_{11} \Theta' = I_J$ . Thus, since  $\Sigma_{11} = D_p$  is diagonal with the class proportions as diagonal elements, we have  $\Theta' \Theta'^T = D_p^{-1}$ , or  $\|\theta^j\|^2 + 1 = 1/p_j \forall j$ :

$$(41) \quad d(h, m^j) = h^T \Sigma_{22}^{-1} h + \|D_{(1-\alpha^2)^{-1/2}} (\theta^j - \eta)\|^2 - 1/p_j + 1.$$

This is the form given by Breiman and Ihaka (1984).

If we neglect terms which are independent of the classes, such as  $h^T \Sigma_{22}^{-1} h$ , the original class-adjusted Mahalanobis distance (26) is seen to be equivalent to the distances labeled 1–5 in Section 3.6.

**Acknowledgments.** We thank Michael Leblanc and Werner Stuetzle for helpful discussion. We also thank the referees and an Associate Editor for their insightful suggestions on an earlier draft of this paper, which encouraged us to rethink some of the issues and resulted in the inclusion of Section 4. This work was done while Trevor Hastie was with AT & T Bell Labs and Andreas Buja was with Bellcore.

## REFERENCES

- BREIMAN, L. and IHAKA, R. (1984). Nonlinear discriminant analysis via scaling and ACE. Technical report, Univ. California, Berkeley.
- BUJA, A., HASTIE, T. and TIBSHIRANI, R. (1989). Linear smoothers and additive models (with discussion). *Ann. Statist.* **17** 453–555.
- CAMPBELL, N. (1980). Shrunk estimators in discriminant and canonical variate analysis. *J. Roy. Statist. Soc. Ser. C* **29** 5–14.
- DE LEEUW, J., YOUNG, F. and TAKANE, Y. (1976). Additive structure in qualitative data: an alternating least squares method with optimal scaling features. *Psychometrika* **41** 471–503.
- DIPILO, P. (1976). The application of bias to discriminant analysis. *Comm. Statist. Theory Methods* **1** 843–854.
- DIPILO, P. (1979). Biased discriminant analysis: evaluation of the optimum probability of classification. *Comm. Statist. Theory Methods* **8** 1447–1458.
- FRIEDMAN, J. (1989). Regularized discriminant analysis. *J. Amer. Statist. Assoc.* **84** 165–175.
- GIFI, A. (1981). Nonlinear multivariate analysis. Unpublished manuscript.
- GIFI, A. (1990). *Nonlinear Multivariate Analysis*. Wiley, New York.
- GREEN, P. and YANDELL, B. (1985). Semi-parametric generalized linear models. *Proceedings of the 2nd International GLIM Conference. Lecture Notes in Statist.* **32** 44–55. Springer, New York.
- HASTIE, T. and MALLOWS, C. (1993). Comment on “A statistical view of some chemometric regression tools” by J. Friedman and I. Frank. *Technometrics* **35** 140–143.
- HASTIE, T., TIBSHIRANI, R. and BUJA, A. (1994). Flexible discriminant analysis by optimal scoring. *J. Amer. Statist. Assoc.* **89** 1255–1270.
- KIIVERI, H. (1992). Canonical variate analysis of high-dimensional spectral data. *Technometrics* **34** 321–331.
- LEBART, L., MORINEAU, A. and WARWICK, K. (1984). *Multivariate Descriptive Statistical Analysis*. Wiley, New York.
- LE CUN, Y., BOSER, B., DENKER, J. S., HENDERSON, D., HOWARD, R., HUBBARD, W. and JACKEL, L. (1990). Handwritten digit recognition with a back-propagation network. In *Advances in Neural Information Processing Systems* (D. Touretzky, ed.) **2**. Morgan Kaufman, Denver.
- LEURGANS, S., MOYEED, R. and SILVERMAN, B. (1993). Canonical correlation analysis when the data are curves. *J. Roy. Statist. Soc. Ser. B* **55** 725–740.
- MARDIA, K., KENT, J. and BIBBY, J. (1979). *Multivariate Analysis*. Academic, New York.
- O’SULLIVAN, F. (1991). Discretized Laplacian smoothing by Fourier methods. *J. Amer. Statist. Assoc.* **86** 634–642.
- RAMSAY, J. and DALZELL, C. (1991). Some tools for functional data analysis (with discussion). *J. Roy. Statist. Soc. Ser. B* **53** 539–572.
- SEBER, G. (1984). *Multivariate Observations*. Wiley, New York.

- TAKANE, Y., YOUNG, F. and DE LEEUW, J. (1979). Nonmetric common factor analysis: an alternating least squares method with optimal scaling features. *Behaviormetrika* **6** 45–56.
- VINOD, H. (1976). Canonical ridge and econometrics of joint production. *J. Econometrics* **4** 147–166.
- VINOD, H. and ULLAH, A. (1981). *Recent Advances in Regression Methods*. Dekker, New York.
- WAHBA, G. (1990). *Spline Models for Observational Data*. SIAM, Philadelphia.
- YOUNG, F., TAKANE, Y. and DE LEEUW, J. (1978). The principal components of mixed measurement level multivariate data: an alternating least squares method with optimal scaling features. *Psychometrika* **43** 505–529.

TREVOR HASTIE  
DEPARTMENT OF STATISTICS  
SEQUOIA HALL  
STANFORD UNIVERSITY  
STANFORD, CALIFORNIA 94305

ANDREAS BUJA  
AT & T BELL LABORATORIES  
ROOM 2C-261  
600 MOUNTAIN AVENUE  
MURRAY HILL, NJ 07974

ROB TIBSHIRAMI  
DEPARTMENT OF PREVENTIVE MEDICINE AND BIostatISTICS  
AND DEPARTMENT OF STATISTICS  
UNIVERSITY OF TORONTO  
TORONTO, ONTARIO  
CANADA