

Flexible Discriminant Analysis by Optimal Scoring

Trevor HASTIE, Robert TIBSHIRANI, AND Andreas BUJA*

Fisher's linear discriminant analysis is a valuable tool for multigroup classification. With a large number of predictors, one can find a reduced number of discriminant coordinate functions that are "optimal" for separating the groups. With two such functions, one can produce a classification map that partitions the reduced space into regions that are identified with group membership, and the decision boundaries are linear. This article is about richer nonlinear classification schemes. Linear discriminant analysis is equivalent to multiresponse linear regression using optimal scorings to represent the groups. In this paper, we obtain nonparametric versions of discriminant analysis by replacing linear regression by any nonparametric regression method. In this way, any multiresponse regression technique (such as MARS or neural networks) can be postprocessed to improve its classification performance.

KEY WORDS: Classification; Discriminant analysis; Nonparametric regression; MARS.

1. INTRODUCTION

Multigroup classification or discrimination is an important problem with applications in many fields. In the generic problem, the outcome of interest G falls into J unordered classes, which for convenience we denote by the set $\mathcal{J} = \{1, 2, 3, \dots, J\}$. We wish to build a rule for predicting the class membership of an item based on p measurements of predictors or features $\mathbf{X} \in R^p$. Our training sample consists of the class membership and predictors for N items. Traditional statistical methods for this problem include linear discriminant analysis and multiple logistic regression. Neural network classifiers have become a powerful alternative, with the ability to incorporate a very large number of features in an adaptive nonlinear model. Ripley (1994) gave an informative survey from a statistician's viewpoint. The recent success and popularity of neural networks led us to look for similar methodologies in the statistical literature, but this seems to be a relatively unexplored area. One significant approach is the classification and regression tree (CART) methodology of Breiman, Friedman, Olshen, and Stone (1984), which is well known to statisticians and is becoming popular in the artificial intelligence community.

There have been a number of recent advances in the nonparametric multiple regression literature. These include projection pursuit regression (Friedman and Stuetzle 1981), the ACE algorithm (Breiman and Friedman 1985), additive models (Hastie and Tibshirani 1990), multivariate adaptive regression splines (MARS; Friedman 1991), Breiman's (1991) Π method, the interaction spline methodology of Wahba (1990), and more recently the hinging hyperplanes of Breiman (1991a). Neural networks (e.g., Barron and Barron 1988; Lippman 1989, and Hinton 1989) can be viewed as yet another approach to nonparametric regression. In this article we describe methods for multigroup classification that use these tools to generalize linear discriminant analysis.

The foundations for the developments described here can be found in the nonlinear scaling literature, notably the work of Gifi (1981, 1990). Our work was motivated by the unpublished paper by Breiman and Ihaka (1984). Section 6.3 details the connection with their work. Ripley and Hjort (1994) were similarly motivated.

This article focuses on adaptive classification procedures. A companion article, "Penalized Discriminant Analysis" (Hastie, Buja, and Tibshirani 1994), gives a more technical basis for some of the procedures described here and focuses on obtaining smooth, interpretable canonical variates for high-dimensional problems such as spectral and image analysis. Both articles rely on the connections between penalized optimal scoring and penalized discriminant analysis. Hereafter we will refer to this companion article as PDA.

2. LINEAR DISCRIMINANT ANALYSIS AND GENERALIZATIONS

2.1 A Review of Linear Discriminant Analysis

Linear discriminant analysis (LDA) is a time-honored, standard tool for classification. A test observation with predictor \mathbf{X}_0 is classified to the class with centroid closest to \mathbf{X}_0 , where distance is measured in the *Mahalanobis* metric using the pooled within-group covariance matrix. This procedure can be justified by assuming that the predictors have a multivariate Gaussian distribution, with different means but a common covariance matrix among the classes. The observation is then assigned to the class having the maximum posterior class probability; this results in the rule described earlier if the class prior probabilities are the same. Variations in the distance thresholds occur if the class prior probabilities are not equal; further generalizations involve using loss functions other than the 0-1 loss implicit here. Two characteristics of this procedure are:

1. All the relevant distance information is contained in the at most $J - 1$ -dimensional subspace of R^p spanned by the J group centroids.
2. The decision boundaries are linear.

A reduced form of LDA due to Fisher and Rao adds a graphical component to the procedure. One finds the $K < J - 1$ dimensional subspace of R^p in which the group centroids

* Trevor Hastie is Staff Member, Statistics and Data Analysis Research Group, AT&T Bell Laboratories, Murray Hill, NJ 07974. Robert Tibshirani is Professor, Department of Preventive Medicine and Biostatistics, and Department of Statistics, University of Toronto, Ontario, Canada M5S 1A8. Andreas Buja is Staff Member, Statistics and Data Analysis Research Group, Bellcore, Morristown, NJ 07960-1910. The authors thank Leo Breiman and Ross Ihaka for a copy of their technical report, and Leo, Jerry Friedman, and Brian Ripley for helpful discussion. They also thank an associate editor and two referees for their helpful and constructive comments. The second author was supported by the Natural Sciences and Engineering Research Council of Canada.

are most separated (once again using the Mahalanobis metric confined to this subspace), and then classifies new data to the closest centroid in the reduced space. This leads to a third characteristic:

3. For small k , the data can be plotted in the reduced space, giving a graphical representation of the group separation.

Often two or three dimensions are all that are needed, even for large J . Finally, in certain problems, typically having many classes and limited training data, the reduced space can be more stable and yield improved misclassification results on test data, demonstrating a fourth characteristic:

4. The dimension-reduced model can show better classification performance.

Figure 1 shows several views of the first 5 LDA variates for 10-dimensional speech data, for which there are 11 classes. The 10 predictors are derived from the spectra of the digitized spoken vowels, which fall into 11 classes. The class separation is evident, especially in the plots of the leading LDA variates. These data are described in more detail in Section 3.4.

The K -dimensional subspace is defined by a $p \times K$ matrix \mathbf{U} of LDA vectors \mathbf{u}_k , and the derived LDA variables are given by $\mathbf{U}^T \mathbf{x}$: similarly the transformed centroids are $\mathbf{U}^T \mathbf{m}_j$ for the j th class. The transformed variables $\mathbf{U}^T \mathbf{x}$ have identity covariance within groups, and the Mahalanobis distance from \mathbf{x} to the j th centroid in this space is simply the Euclidean distance $\delta(\mathbf{x}, \mathbf{m}_j) = \|\mathbf{U}^T(\mathbf{x} - \mathbf{m}_j)\|$. In words, we first “sphere” the data using the common within-groups covariance matrix, project these data onto the $J - 1$ -dimensional subspace spanned by the J (sphered) centroids (or a subspace thereof), and then classify new observations to the class corresponding to the closest centroid.

In practice the centroids \mathbf{m}_j and within-groups covariance matrix Σ_W are unknown and must be estimated from the data. The LDA vectors \mathbf{u}_k successively diagonalize the quadratic form $\mathbf{u}^T \Sigma_{\text{Bet}} \mathbf{u}$ with normalization constraint $\mathbf{u}^T \Sigma_W \mathbf{u} = 1$. Here $\Sigma_{\text{Bet}} = \sum_j (N_j/N)(\mathbf{m}_j - \bar{\mathbf{m}})(\mathbf{m}_j - \bar{\mathbf{m}})^T$ is the *between-groups* covariance matrix. Gnanadesikan and Kettenring (1989) gave a recent review of the current literature on discriminant analysis.

2.2 Classification by Multivariate Linear Regression

Multiresponse linear regression can also be used for classification. We form an indicator response matrix having J columns, with a value 1 in the j th column if the observation is in class j and 0 otherwise. Then we carry out a multivariate linear regression of the response matrix on the predictors and classify each observation by the class having the largest fitted value $\hat{Y}_1, \hat{Y}_2, \dots, \hat{Y}_J$. Although \hat{Y}_j can be thought of as an estimate of the class- j probability, there is no guarantee that it lies in $[0, 1]$, nor in general that the J estimates sum to 1. For this reason a modified estimate $\hat{p}_j = \exp(\hat{Y}_j) / \sum_i \exp(\hat{Y}_i)$ is often preferred. In the neural network literature, this is known as *softmax* (Bridle 1990), and we will

use that terminology here. In the two-group case, with equal sample sizes, *softmax* is essentially equivalent to LDA. They are not equivalent in general, and *softmax* does not seem to work as well as linear discriminant analysis in the examples we have tried. For example, in the vowel example of Figure 1, LDA had a test error rate of 56%, whereas *softmax* had a test error rate of 67%. In Section 5 we shed some light on the reason for this difference.

We prove in the Appendix that the space of linear discriminant fits is the same as the space of multivariate linear regression fits, as defined previously. This means that the LDA solution can be obtained from a linear discriminant analysis of the *fitted values* from the multivariate regression fits. This fact is clearly known to some (Breiman and Ihaka [1984] gave a proof, and we discuss the equivalences in detail in PDA), but does not seem to be widely appreciated. Using LDA in this fashion as a postprocessor for multivariate linear regression generally improves its classification performance, sometimes dramatically; see Section 5 for a detailed discussion of this point.

2.3 Generalization of Linear Discriminant Analysis

In practice linear decision boundaries are often too crude, and nonlinear boundaries can be more effective. The Gaussian assumptions are rarely met, and sometimes a group might even be disjoint. Using different class-covariance matrices in the Bayesian procedure results in quadratic decision boundaries. A different approach is to augment the predictor set to include quadratic and bilinear terms and then perform the LDA in the enlarged space; this also leads to quadratic decision boundaries in the original space, although slightly more restrictive than in the former. Besides the fact that a globally quadratic boundary may not adapt well to the problem at hand, the generalization from linear to quadratic discrimination adds $O(Jp^2)$ parameters to the model, far too many if p is large.

We propose to use nonparametric regression procedures to estimate nonlinear boundaries for classification. To this end, we make use of the well-known fact that LDA is equivalent to canonical correlation analysis; the linear predictors define the one set of variables, and a set of dummy variables representing class membership defines the other set. Canonical correlation analysis in this context gives the solution to a scoring problem that we now describe.

Suppose that $\theta: \mathcal{J} \mapsto R^1$ is a function that assigns scores to the classes such that the transformed class labels are optimally predicted by linear regression on \mathbf{X} . This produces a one-dimensional separation between the classes. More generally, we can find K sets of independent scorings for the class labels, $\{\theta_1, \theta_2, \dots, \theta_K\}$ and K corresponding linear maps $\eta_k(\mathbf{X}) = \mathbf{X}^T \beta_k$, $k = 1, \dots, K$, chosen to be optimal for multiple regression in R^K . If our training sample has the form (g_i, \mathbf{x}_i) , $i = 1, 2, \dots, N$ then the scores $\theta_k(g)$ and the maps β_k are chosen to minimize the average squared residual:

$$ASR = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N (\theta_k(g_i) - \mathbf{x}_i^T \beta_k)^2. \quad (1)$$

The scores are assumed to be mutually orthogonal and normalized with respect to an appropriate inner product to pre-

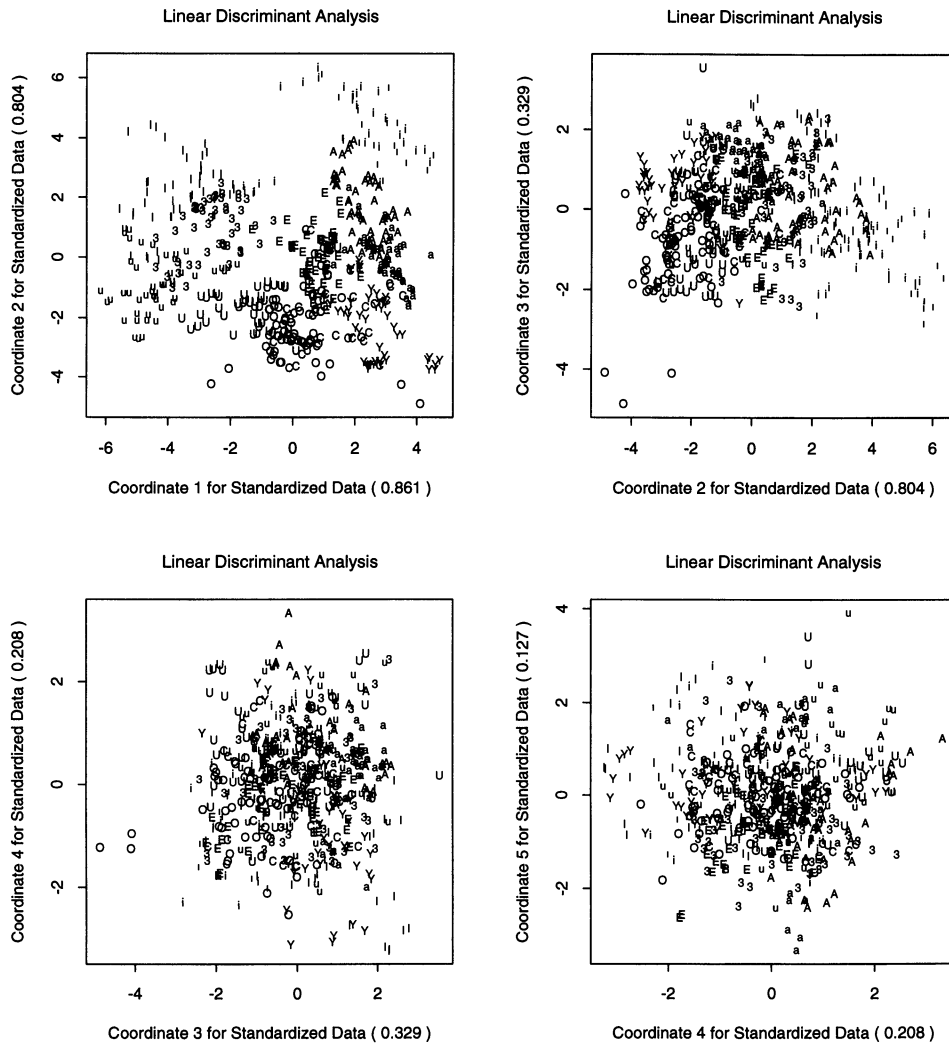


Figure 1. Some Selected Two-Dimensional Projections of the LDA Subspace Computed for the Vowel Data Described in Section 3.4. The eigenvalues for each coordinate are in parentheses, and indicate the strength of that variable in separating the groups. These plots show visually that most of the discrimination lies in the first few discriminant variables.

vent trivial zero solutions. The key fact that we use is the following:

The sequence of LDA vectors \mathbf{u}_k is identical to the sequence β_k up to a constant (Mardia, Kent, and Bibby 1979).

The standard way of carrying out a canonical correlation analysis is by way of a suitable singular value decomposition. Here we outline a somewhat different approach to the solution that lends itself naturally to a nonlinear generalization.

Let \mathbf{Y} be the $N \times J$ indicator matrix corresponding to the dummy-variable coding for the classes; that is, the ij th element of \mathbf{Y} is 1 if the i th observation falls in class j , and 0 otherwise. Let $\Theta_{J \times K}$, $K \leq J - 1$ be a matrix of K score vectors for the J classes. If we let Θ^* be the $N \times K$ matrix of transformed values of the classes with ik th element $\theta_k(g_i)$, then $\Theta^* = \mathbf{Y}\Theta$. Looking at (1), it is clear that if the scores were fixed, we could minimize ASR by regressing Θ^* on \mathbf{x} . If we let \mathbf{P}_X project onto the column space of the predictors, this says that

$$ASR(\Theta) = \text{tr}\{\Theta^{*T}(I - \mathbf{P}_X)\Theta^*\} / N \\ = \text{tr}\{\Theta^T \mathbf{Y}^T (I - \mathbf{P}_X) \mathbf{Y} \Theta\} / N. \quad (2)$$

If we assume that the scores have mean zero, unit variance, and are uncorrelated for the N observations ($\Theta^{*T}\Theta^* / N = \mathbf{I}_K$), minimizing (2) amounts to finding the K largest eigenvectors Θ of $\mathbf{Y}^T \mathbf{P}_X \mathbf{Y}$ with normalization $\Theta^T \mathbf{D}_p \Theta = \mathbf{I}_K$, where $\mathbf{D}_p = \mathbf{Y}^T \mathbf{Y} / N$, a diagonal matrix of the sample class proportions N_j / N . We could do this by constructing the matrix \mathbf{P}_X , computing $\mathbf{Y}^T \mathbf{P}_X \mathbf{Y}$, and then calculating its eigenvectors. But a more convenient approach avoids explicit construction of \mathbf{P}_X and takes advantage of the fact that \mathbf{P}_X computes a linear regression. Here is a summary:

Linear discriminant analysis by optimal scoring: LDA

1. **Initialize.** Form \mathbf{Y} , the $N \times J$ indicator matrix corresponding to the dummy-variable coding for the classes.
2. **Multivariate regression.** Set $\hat{\mathbf{Y}} = \mathbf{P}_X \mathbf{Y}$ and denote the $p \times J$ coefficient matrix by \mathbf{B} : $\hat{\mathbf{Y}} = \mathbf{X}\mathbf{B}$.
3. **Optimal scores.** Obtain the eigenvector matrix Θ of $\mathbf{Y}^T \hat{\mathbf{Y}} = \mathbf{Y}^T \mathbf{P}_X \mathbf{Y}$ with normalization $\Theta^T \mathbf{D}_p \Theta = \mathbf{I}$.

4. **Update.** The coefficient matrix in step 2 to reflect the optimal scores: $\mathbf{B} \leftarrow \mathbf{B}\mathbf{\Theta}$. The final optimally scaled regression fit is the $(J - 1)$ vector function $\eta(\mathbf{x}) = \mathbf{B}^T \mathbf{x}$.

This is an alternative algorithm for computing the usual canonical variates. The final coefficient matrix \mathbf{B} is, up to a diagonal scale matrix, the same as the discriminant analysis coefficient matrix. Specifically, $\mathbf{U}^T \mathbf{x} = \mathbf{D}\mathbf{B}^T \mathbf{x} = \mathbf{D}\eta(\mathbf{x})$, where

$$D_{kk}^2 = \frac{1}{[\alpha_k^2(1 - \alpha_k^2)]} \quad (3)$$

and α_k^2 is the k th largest eigenvalue computed in step 3. A derivation is given in the Appendix.

Our interest in this procedure is that it provides a starting point for generalizing LDA to a nonparametric version; we replace the linear-projection operator \mathbf{P}_X by a nonparametric regression procedure, which we denote by the linear operator S . One simple and effective approach toward this end is to expand \mathbf{X} into a larger set of basis variables $h(\mathbf{X})$ and then simply use $\mathbf{S} = P_{h(\mathbf{X})}$ in place of \mathbf{P}_X . The second approach to quadratic discriminant analysis is a member of this class, in which case h is a second-degree polynomial expansion.

The procedures that we have in mind are more “nonparametric” than this. In this article we concentrate on two strategies for adaptive nonparametric regression: MARS and BRUTO, both described in detail in Section 4. MARS is an adaptive regression technique able to capture interactions in a hierarchical manner (Friedman 1991). BRUTO is an adaptive method for estimating an additive model of the form $\alpha + \sum_1^p f_j(X_j)$ using smoothing splines. In terms of basis functions, they can be characterized as follows:

- MARS builds up a basis expansion $h(\mathbf{X})$ adaptively, including particular localized interactions only where needed.
- BRUTO automatically generates a very large basis set and achieves parsimony by shrinking coefficients in a judicious and structured way.

In both of these procedures, if one conditions on the adaptive choices of terms and/or smoothing parameters, then the fitting mechanisms can be represented as the action of a linear operator acting on the response. This is a necessary property for their use here. We denote these choices/smoothing parameters collectively by the *hypersmoothing* parameter λ and, given λ , denote the linear smoothing operator by $S(\lambda)$.

The real payoff in expressing the discriminant analysis model in a regression context is the ability to perform model selection and regularization. MARS and BRUTO belong to an ever-growing class of adventurous and adaptive regression procedures (including neural networks), which we can simply “plug into” the discriminant analysis problem and inherit all the adaptivity in the new context.

Before we give details specific to the particular methods, we present our general proposal.

Flexible discriminant analysis by optimal scoring: FDA

1. **Initialize.** Choose an initial score matrix $\mathbf{\Theta}_0$ with $K \leq J$ satisfying the constraints $\mathbf{\Theta}^T \mathbf{D}_p \mathbf{\Theta} = \mathbf{I}$; let $\mathbf{\Theta}_0^* = \mathbf{Y}\mathbf{\Theta}_0$.

2. **Multivariate nonparametric regression.** Fit a multi-response adaptive nonparametric regression of $\mathbf{\Theta}_0^*$ on \mathbf{X} , giving fitted values $\hat{\mathbf{\Theta}}_0^*$. Let $S(\hat{\lambda})$ be the linear operator that fits the final chosen model and let $\eta(\mathbf{x})$ be the vector of fitted regression functions.

3. **Optimal scores.** Obtain the eigenvector matrix Φ of $\mathbf{\Theta}_0^{*T} \hat{\mathbf{\Theta}}_0^* = \mathbf{\Theta}_0^{*T} S(\hat{\lambda}) \mathbf{\Theta}_0^*$, and hence the optimal scores $\mathbf{\Theta} = \mathbf{\Theta}_0 \Phi$.

4. **Update** the final model from step 2 using the optimal scores: $\mathbf{n}(\mathbf{x}) \leftarrow \Phi^T \eta(\mathbf{x})$.

For a J class problem, the vector of canonical variates or functions $\eta(\mathbf{x})$ has at most $K = J - 1$ components. If $\bar{\eta}^j = \sum_{g_i=j} \eta(\mathbf{x}_i) / N_j$ denotes the fitted centroid of the j th class in this space of canonical variates, then the discrimination rule has the form of a (weighted) nearest centroid rule:

Assign an observation x to the class j that minimizes

$$\delta(\mathbf{x}, \mathbf{j}) = \|\mathbf{D}(\eta(\mathbf{x}) - \bar{\eta}^j)\|^2. \quad (4)$$

\mathbf{D} is once again the diagonal matrix (3) of scale factors that convert optimally scaled fits to discriminant analysis variables.

We have glossed over many details, which we now discuss.

2.4 More Details on $S(\lambda)$

$S(\lambda)$ is a linear operator (in this case an $N \times N$ matrix), and λ is the hyperparameter that captures the model selection/smoothing parameter selection required to specify the amount and type of smoothing. $S(\lambda)\mathbf{\Theta}^*$ for a $N \times K$ matrix of “responses” $\mathbf{\Theta}^*$ implies matrix multiplication. In practice we have efficient algorithms for computing $S(\lambda)\mathbf{\Theta}^*$ without explicitly computing $S(\lambda)$ itself.

For example, in an additive spline model with p terms, λ_j might be the roughness penalty for the j th term. The “hat” in $S(\hat{\lambda})$ indicates that we have used adaptive techniques to select the values for λ .

Note that the foregoing description implies that the *same procedure* $S(\lambda)$ is used to fit each of the K models (one for each set of scorings); in particular, the same value for λ is used. Therefore, the nonparametric regression procedure must be able to handle a multiple response variable when selecting λ . We discuss the relevant details in the sections on MARS and BRUTO.

Although we have chosen to emphasize MARS and BRUTO in this article, there are many other possibilities for $S(\lambda)$, including projection pursuit, neural networks, and hinging hyperplanes. The theory works for linear operators $S(\lambda)$, and none of the other approaches is strictly linear. Selecting $\hat{\lambda}$ from the training data typically makes $S(\hat{\lambda})$ non-linear (because $\hat{\lambda}$ is a function of \mathbf{Y}). Our theory applies only if we condition on these choices. It does offer some comfort, however, that the technique might be useful as a postprocessor for *any* multivariate regression technique.

2.5 The Canonical Functions η

The nonparametric procedures produce canonical functions $\eta(\mathbf{x})$, whereas in the linear case LDA produces discriminant variates $\mathbf{U}^T \mathbf{x}$. These latter variates are scaled versions of the *linear* canonical functions, $\mathbf{U}^T \mathbf{x} = \mathbf{D}\mathbf{B}^T \mathbf{x}$, and

hence $\mathbf{D}\eta(\mathbf{x})$ are nonparametric versions of the discriminant variates. (We avoid calling them “discriminant functions,” a phrase reserved for distance functions.)

For both nonparametric procedures that we use, the functions $\eta(\mathbf{x})$ are in fact linear combinations of an adaptively chosen set of basis functions $h(\mathbf{x})$:

$$\eta(\mathbf{x}) = \mathbf{B}^T h(\mathbf{x}).$$

This suggests that given the basis functions (selected in step 2), we can simply perform LDA, treating them as variables. When the MARS algorithm or any projection is used for $\mathbf{S}(\lambda)$, this is exactly the case; the solutions resulting from step 4 are exactly a LDA using the basis variables $h(\mathbf{x})$. Similarly, we could instead perform LDA using the fitted values \hat{Y}_j from the MARS or projection fit as variables (because $P_{\hat{Y}}\mathbf{Y} = P_{h(\mathbf{X})}\mathbf{Y}$.) Neither of these equivalences apply when regularized (shrunk) regression procedures such as BRUTO are used; in these cases the FDA procedure solves a form of *penalized discriminant analysis* (PDA) described next.

2.6 Penalized Discriminant Analysis

The idea here is to first expand the predictors into a (much) larger set of basis functions $h(\mathbf{x})$, but then ensure by regularization that the linear combinations $\beta^T h(\mathbf{x})$ are in some sense smooth. To focus ideas, we use the additive spline model produced by the BRUTO procedure.

Suppose that we use an additive model $\sum f_j(X_j)$ for the scores in step 2 of the FDA procedure. To formalize the estimation of f_j 's, we consider a penalized version of the averaged squared residual (1):

$$\frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N \left[\theta_k(g_i) - \sum_{j=1}^p f_{jk}(x_{ij}) \right]^2 + \sum_{k=1}^K \sum_{j=1}^p \lambda_j \int f_{jk}''(t)^2 dt, \quad (5)$$

where x_{ij} is the measurement on the j th variable for the i th observation and λ_j are smoothing parameters that govern the trade-off between fit and smoothness. The solution is a finite-dimensional additive-spline model, and (5) can be written as

$$\frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N \left[\theta_k(g_i) - \sum_{j=1}^p \beta_{jk}^T \mathbf{h}_j(x_{ij}) \right]^2 + \sum_{k=1}^K \sum_{j=1}^p \lambda_j \beta_{jk}^T \boldsymbol{\Omega}_j \beta_{jk}, \quad (6)$$

where each \mathbf{h}_j a vector of up to N natural-spline basis functions defined on the set $\{x_{1j}, \dots, x_{Nj}\}$ and β_{jk} , $k = 1, \dots, K$ are the corresponding coefficients. Thus the entire basis set $h(\mathbf{X})$ can have as many as Np components, in which case the ultimate block-diagonal penalty matrix $\boldsymbol{\Omega}(\lambda) = \text{diag}(\lambda_1 \boldsymbol{\Omega}_1, \dots, \lambda_p \boldsymbol{\Omega}_p)$ will have dimension $Np \times Np$ and the coefficient matrix \mathbf{B} dimension $Np \times K$. In this case the regression operator has the form

$$\mathbf{S}(\lambda) = \mathbf{H}(\mathbf{H}^T \mathbf{H} + \boldsymbol{\Omega}(\lambda))^{-1} \mathbf{H}^T,$$

and the partially optimized criterion (6) reduces to

$$ASR_p(\boldsymbol{\Theta}) = \text{tr}\{\boldsymbol{\Theta}^T \mathbf{Y}^T (\mathbf{I} - \mathbf{S}(\lambda)) \mathbf{Y} \boldsymbol{\Theta}\} / N. \quad (7)$$

The BRUTO procedure described in Section 4.2 is an efficient algorithm for computing the fit $\mathbf{S}(\lambda) \boldsymbol{\Theta}^*$ (for fixed θ_k) and selecting the smoothing parameters λ_j adaptively.

The problem solved by minimizing Equation (6) corresponds to a form of *penalized discriminant analysis*. Let $\boldsymbol{\Sigma}_{\text{Bet}}$ be the between-groups covariance matrix for $h(\mathbf{X})$ and let $\boldsymbol{\Sigma}_w + \boldsymbol{\Omega}$ be the *penalized* within-groups covariance. We define:

A *penalized discriminant analysis* finds a matrix \mathbf{U} to maximize $\text{tr}(\mathbf{U}^T \boldsymbol{\Sigma}_{\text{Bet}} \mathbf{U})$ subject to the penalized constraint $\mathbf{U}^T (\boldsymbol{\Sigma}_w + \boldsymbol{\Omega}) \mathbf{U} = \mathbf{I}$.

The discrimination rule is given by (4). Without $\boldsymbol{\Omega}$ this procedure is just LDA on $h(\mathbf{X})$. The inclusion of $\boldsymbol{\Omega}$ forces smoothness of the resulting solution and prevents overfitting. The procedure can also be interpreted in terms of penalized Mahalanobis distances from class centroids in the augmented space of $h(\mathbf{X})$:

$$\begin{aligned} \delta(\mathbf{x}, j) &= (h(\mathbf{x}) - \bar{h}^j)^T (\boldsymbol{\Sigma}_w + \boldsymbol{\Omega})^{-1} (h(\mathbf{x}) - \bar{h}^j) \\ &= \|\mathbf{D}(\eta(\mathbf{x}) - \bar{\eta}^j)\|^2, \end{aligned} \quad (8)$$

where (8) is the same as (4).

An aside: In some applications we have basis functions h that represent the actual input features, rather than being an expanded set of predictors derived from the input features. For example, the h might be the grayscale level of a pixel in an image. In this case we wish the *coefficients* themselves to form a smooth image (neighboring pixels carry similar and highly correlated information.) A procedure similar to that described earlier leads to a ridge-regression type estimate for the β_j s. This is the central topic of the PDA paper and will not be discussed further here.

2.7 Initial Values for $\boldsymbol{\Theta}$

In the LDA algorithm on page 10 there was no initialization, whereas for the nonparametric version we have used an initial $\boldsymbol{\Theta}_0$. We first summarize the issues before we get into details:

- $\boldsymbol{\Theta}_0 = \mathbf{I}_J$ is a valid choice and thus amounts to using $\boldsymbol{\Theta}_0^* = \mathbf{Y}$ as in LDA.
- Any rank $(J-1)\boldsymbol{\Theta}_0$ of dimension $J \times (J-1)$ which satisfies the normalization constraints and is orthogonal to $\mathbf{1}_J$ (a J -vector of 1s), is equivalent to $\boldsymbol{\Theta}_0 = \mathbf{I}$. The only advantage is that the multivariate smoother $\mathbf{S}(\lambda)$ has one less column to fit.
- If $\boldsymbol{\Theta}_0$ has $K < J-1$ column and we seek a K -dimensional solution, then starting values can be critical, as they can alter the search for λ .

The first claim is obvious. To verify the second claim, we first note that $\mathbf{Y} \mathbf{1}_J = \mathbf{1}_N$ and for most sensible regression procedures, $\mathbf{S}(\lambda) \mathbf{1}_N = \mathbf{1}_N$ for any λ . This also implies that $\mathbf{1}_J$ is a trivial eigenvector of $\mathbf{Y}^T \mathbf{S}(\lambda) \mathbf{Y}$. The eigendecomposition is equivariant under a rotation $\mathbf{R} = (\mathbf{1}_J; \boldsymbol{\Theta}_0)$, which maps \mathbf{Y} to $\mathbf{Y} \mathbf{R} = (\mathbf{1}_N; \mathbf{Y} \boldsymbol{\Theta}_0) = (\mathbf{1}_N; \boldsymbol{\Theta}_0^*)$. It is expedient to have the $J \times (J-1)$ contrast matrix $\boldsymbol{\Theta}_0$ satisfy the normalization conditions required for the discriminant analysis: $\boldsymbol{\Theta}_0^T \mathbf{1}_J = 0$, $\boldsymbol{\Theta}_0^T \mathbf{D}_p \boldsymbol{\Theta}_0 = \mathbf{I}_{J-1}$. The eigenvalue problem thus re-

duces to decomposing the $(J - 1) \times (J - 1)$ matrix $\Theta_0^*{}^T S(\hat{\lambda}) \Theta_0^*$.

All that remains is to show that Θ_0 does not affect the selection of λ . Given Θ_0 , the *ASR* criterion is

$$\begin{aligned} ASR &= \frac{1}{N} \text{tr} \{ \Theta_0^T Y^T (I - S(\lambda))^T (I - S(\lambda)) Y \Theta_0 \} \\ &= \frac{1}{N} \text{tr} \{ Y^T (I - S(\lambda))^T (I - S(\lambda)) Y \Theta_0 \Theta_0^T \}. \quad (9) \end{aligned}$$

As we will see, *ASR* is the numerator of the *GCV* criterion used to select λ . The normalization on Θ_0 implies that $\Theta_0 \Theta_0^T = D_p^{-1} - 11^T$ in (9), so any legitimate version of Θ_0 will have the same effect on (9). Another interpretation is that any rank $J - 1$ legitimate initialization amounts to weighting the j th column of Y by $\sqrt{N/N_j}$ and using a weighted *ASR* criterion in these original dummy variables. This completes the verification of the second claim.

If we seek a $K < J - 1$ -dimensional solution, then an initial $J \times K$ Θ_0 could make a difference, because it could steer the search for λ in important directions. This suggests iterating the entire procedure in some form to improve the choices of Θ_0 . We discuss iteration further in Section 6.1.

2.8 Model Selection: Choosing λ

Nonparametric regression procedures are typically adaptive and require a criterion for model selection; in our terminology we need to choose a value for the hyperparameter λ . We use the generalized cross-validation (*GCV*) criterion

$$GCV(c, \lambda) = \frac{ASR(\lambda)}{[1 - \{1 + c \cdot df(\lambda)\}/N]^2}. \quad (10)$$

Here $df(\lambda)$ is the effective degrees of freedom in the model, not counting the constant term. For the MARS model, $df(\lambda)$ is the number of independent basis functions, whereas for BRUTO it measures the amount of smoothing. In both cases we use $df(\lambda) = \text{tr } S(\lambda) - 1$.

We have introduced an additional hyperparameter c that represents the *cost* per degree of freedom. This will typically be greater than 1 to account for the adaptive nature of the fitting in both cases.

The rationale is as follows. The *GCV* criterion is a regression-based criterion, but our real goal is discrimination or classification. Both MARS and BRUTO (and many other nonparametric regression procedures) are geared-up to select λ via *GCV*, so it is expedient to use *GCV* here as well. By separating out a single cost parameter, c , we can use a regression-based *GCV* to select λ given c and then use a more goal-oriented criterion for selecting c . This is also the basis of the *cost-complexity pruning* of Breiman et al. (1984).

There are several approaches to selecting c . The simplest is to use a fixed value, and, based on the work of Friedman (1991) and Owen (1991) and our experience here, it seems that reasonable values are 2 for additive models (BRUTO and degree-1 MARS models) and 3 for higher-degree MARS models. A more ambitious approach is to choose c to minimize the misclassification cost estimated by m -fold cross-validation. If there are ample data, such sample reuse tech-

niques are not needed, and we can set aside a portion of the data for tuning the procedure.

3. EXAMPLES

In the following examples we compare the new procedures to a number of existing classification methods. The methods include:

1. Linear discriminant analysis (LDA). The classical approach using all the variables.
2. Quadratic discriminant analysis (QDA): Observations are classified to the class with the closest centroid, using Mahalanobis distance based on the class-specific covariance matrix.
3. CART, the classification and regression tree procedure of Breiman et al. (1984). Coordinate splits and default input parameter values are used. In some cases we also report the results using CART with linear combination splits.
4. FDA/BRUTO, flexible discriminant analysis using additive models with adaptive selection of terms and spline smoothing parameters. This is one of the procedures proposed in this article.
5. FDA/MARS, flexible discriminant analysis using Friedman's (1991) multivariate adaptive regression splines. This is another proposal of this article, with an adaptation of the MARS algorithm to allow for multiple response. We report the *degree* of interaction used; *degree* = 1 implies an additive model.

Softmax refers to use of the maximum fitted value for classification, rather than discriminant analysis. *Softmax* is not reported in two-group examples because it agrees with LDA if the class sizes are equal and they are generally in close agreement.

In some cases, neural network procedures are also compared; details are given in the specific example.

3.1 Example: Spherical Clusters

The first of these simulated examples has ten predictors and two classes. The last six predictors are noise variables, with standard normal distributions independent of each other and the class membership. The first four predictors in class 1 are independent standard normal, conditioned on the radius being greater than 3, whereas the first four predictors in class 2 are independent standard normal without the restriction. The first class almost completely surrounds the second class in the four-dimensional subspace of the first four predictors. We chose 250 observations from each class. The error rates of various procedures are shown in Table 1. The optimal decision boundary is the hull of a sphere with radius 3 in the first four coordinates and has an error rate $P(\chi_4^2 > 9) = .061099$.

In this case an additive model does the job exactly, because quadratic coordinate functions in the first four variables suffice to describe the spherical hull. Both the BRUTO and MARS (degree = 1) procedure are additive, they both produced terms involving only variables 1–4, the coordinate functions were roughly quadratic, and they both yielded the smallest misclassification error. The left plot in Figure 2

Table 1. Spherical Distributions in Ten-Dimensional Space With Two Classes

Technique	Error rates	
	Training	Test
LDA	.439 (.006)	.502 (.006)
QDA	.086 (.003)	.138 (.004)
CART	.056 (.002)	.056 (.002)
FDA/BRUTO	.051 (.004)	.061 (.003)
FDA/MARS (degree = 1)	.052 (.003)	.065 (.004)
FDA/MARS (degree = 2)	.042 (.003)	.078 (.005)

NOTE: Six of the variables are noise, whereas the other four define inner and outer spherical regions that define the classes. The results reported below the horizontal line refer to our FDA routine using the named nonparametric regression procedure. The values are averages over 10 simulations, with the standard error of the average in parentheses.

shows the decision region based on the BRUTO fit for the section of ten-space defined by setting variables 3–10 identically equal to zero.

The right plot shows a decision boundary for a similar but simpler classification problem. Equal-sized samples were generated from bivariate Gaussian distributions, each with mean zero and diagonal covariance matrices; the one group has variances four times that of the other. Bayes classification theory tells us the optimal decision boundary is a circle; the figure shows the optimal boundary and indicates the classification found using the BRUTO procedure.

3.2 Example: A Pure Interaction

Here we sampled 200 bivariate uniform variates from $[-1, 1] \times [-1, 1]$. Observations in the southeast and northwest corners were assigned class 1; the others were assigned class 2. The results are shown in Table 2.

If CART chooses its first split along $x_1 = 0$ or $x_2 = 0$, it can achieve perfect prediction after its secondary splits. But

Table 2. Pure Interaction Example Using a Sample of 200 Observations From a Uniform Distribution on $[-1, 1] \times [-1, 1]$, With the Two Classes Defined by the Diagonally Opposite Pairs of Quadrants

Technique	Error rates	
	Training	Test
LDA	.471 (.010)	.485 (.013)
QDA	.029 (.004)	.038 (.005)
CART	.015 (.004)	.497 (.010)
FDA/BRUTO	.500 (.000)	.500 (.000)
FDA/MARS (degree = 1)	.491 (.009)	.500 (.000)
FDA/MARS (degree = 2)	.034 (.004)	.050 (.007)

NOTE: The values are averages over 10 simulations, with the standard error of the average in parentheses.

in its initial split it is unaware of the nested structure and hence does not know to split at $x_1 = 0$ or $x_2 = 0$. The FDA procedure using an additive MARS regression (degree = 1) performed poorly, but it did very well when interactions were allowed (degree = 2).

3.3 Example: Waveform Data

This example is from Breiman et al. (1984, p. 49–55); a concise description was given by Loh and Vanichsetakul (1988). It is a three-class problem with 21 variables. The predictors are defined by

$$x_i = uh_1(i) + (1 - u)h_2(i) + \varepsilon_i \quad \text{Class 1,}$$

$$x_i = uh_1(i) + (1 - u)h_3(i) + \varepsilon_i \quad \text{Class 2,}$$

and

$$x_i = uh_2(i) + (1 - u)h_3(i) + \varepsilon_i \quad \text{Class 3,} \quad (11)$$

where $i = 1, 2, \dots, 21$, u is uniform on $(0, 1)$, ε_i are standard normal variates, and the h_i are the shifted triangular wave-

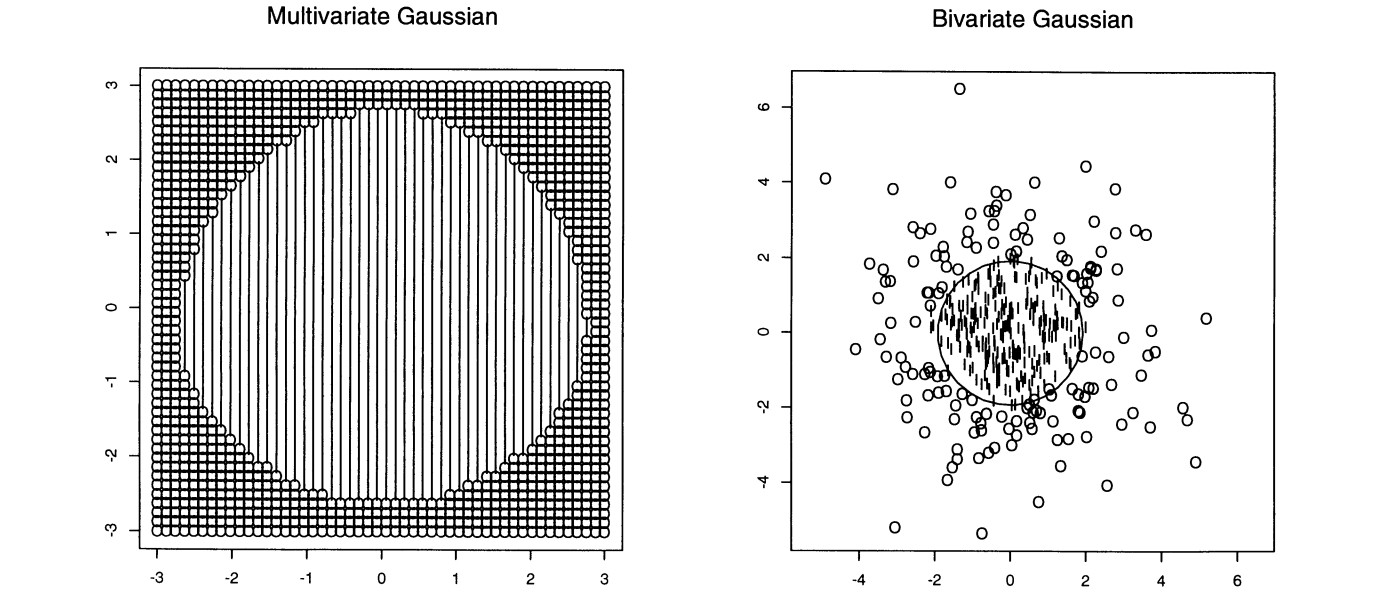


Figure 2. The Left Plot Shows a View of the Decision Boundary for the Ten-Dimensional Gaussian Example Fit Using the BRUTO Procedure, Obtained by Setting Coordinates 3–10 to Zero. The plot symbols are “I” or “O” for inner and outer, produced over a uniform grid in this two-dimensional slice. The right plot shows the predicted class and the Bayes-optimal decision boundary for the two-dimensional Gaussian example.

Table 3. Results for Waveform Data

Technique	Error rates	
	Training	Test
LDA	.121 (.006)	.191 (.006)
Softmax	.121 (.005)	.194 (.007)
QDA	.039 (.004)	.205 (.006)
CART	.072 (.003)	.289 (.004)
FDA/MARS (degree = 1)	.100 (.006)	.191 (.006)
Softmax	.100 (.006)	.191 (.006)
FDA/MARS (degree = 2)	.068 (.004)	.215 (.002)
Softmax	.071 (.003)	.216 (.002)

NOTE: The values are averages over 10 simulations, with the standard error of the average in parentheses.

forms: $h_1(i) = \max(6 - |i - 11|, 0)$, $h_2(i) = h_1(i - 4)$ and $h_3(i) = h_1(i + 4)$.

The training sample has 300 observations, and equal priors were used, so there are roughly 100 observations in each class. We used a test sample of size 500. The results in Table 3 shows that none of the nonlinear methods significantly outperform LDA in terms of misclassification error.

Figure 3 shows the two canonical variates produced in the case of LDA and FDA/BRUTO, evaluated at the test data; the classes are coded in the plot. The LDA plot suggests that the classes lie on the edges of a triangle. This makes sense, because the $h_j(i)$ are represented by three points in 21-space, thereby forming vertices of a triangle. Each class is represented as a convex combination of a pair of vertices and hence lie on an edge. The clusters in the FDA plot are more rounded than in the former; nonlinear transformations permit the bending of the edges of the triangle, thereby achieving slightly better separation.

The results for CART are in qualitative agreement with those reported by Breiman et al. (1984)—test sample error

Table 4. Words Used in Recording the Vowels

Vowel	Word
i	heed
I	hid
E	head
A	had
a:	hard
Y	hud
O	hod
C:	hoard
U	hood
u:	who'd
3:	heard

rates of .28 and .20 for the univariate-split and linear-combination-split versions. They did not report results for full LDA but reported a test error of .24 for stepwise LDA.

3.4 Example: Vowel Recognition Data

This example, a popular benchmark for neural network algorithms, consists of training and test data with 10 predictors and 11 classes. We obtained the data from the benchmark collection maintained by Scott Fahlman at Carnegie Mellon University. The data were contributed by Anthony Robinson (see Robinson 1989), who provided the following description.

An ASCII approximation to the International Phonetic Association symbol and the word in which the 11 vowel sounds were recorded is given in Table 4. The word was uttered once by each of the 15 speakers. Four male and four female speakers were used to train the networks, and the other four male and three female speakers were used for testing the performance.

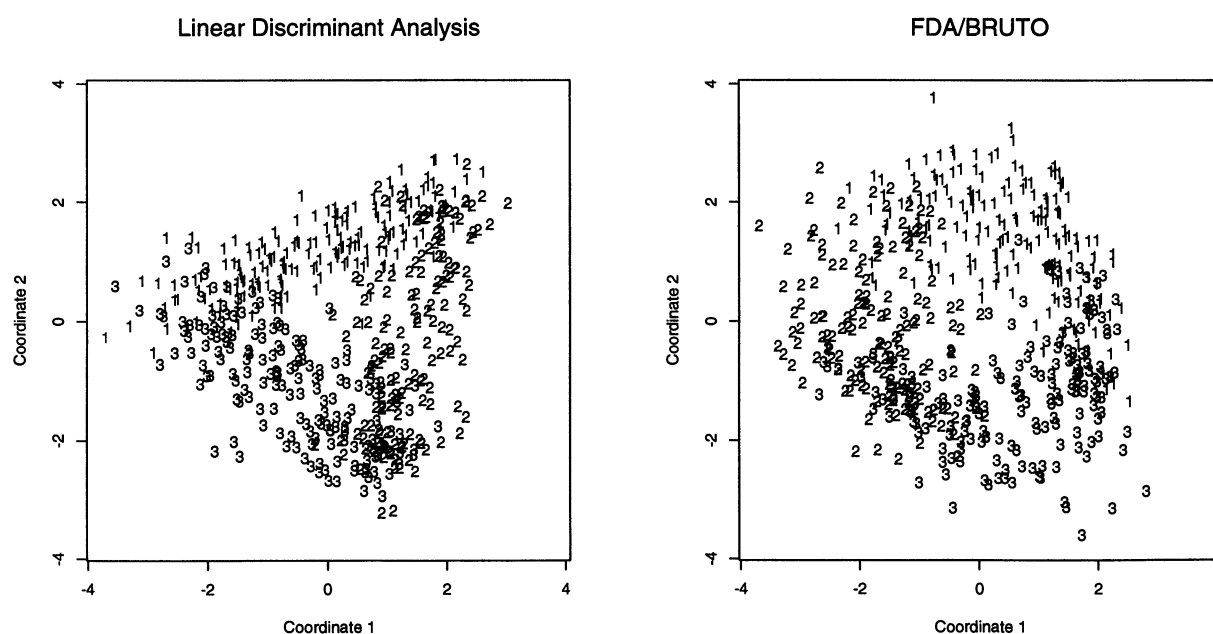


Figure 3. The Left Plot Shows the Canonical Variates Derived From LDA for the Waveform Data, Evaluated for the Test Data. The classes are numerically coded. The right plot is the same, but is derived from the FDA/BRUTO model.

This paragraph is technical and describes how the analog speech signals were transformed into a 10-dimensional feature vector. The speech signals were low-pass filtered at 4.7 kHz and then digitized to 12 bits with a 10 kHz sampling rate. Twelfth-order linear predictive analysis was carried out on six 512-sample Hamming windowed segments from the steady part of the vowel. The reflection coefficients were used to calculate 10 log-area parameters, giving a 10-dimensional input space. (For a general introduction to speech processing and an explanation of this technique, see Rabiner and Schafer 1978.) Each speaker thus yielded six frames of speech from 11 vowels. This gave 528 frames from the eight speakers used to train the networks and 462 frames from the seven speakers used to test the networks.

The results of several classification procedures are shown in Table 5. Lines 5–8 were taken from Robinson (1989); the first three performed the best among the neural network classifiers reported by Robinson. Lines 9–11 show the results of various FDA fits, and line 12 shows the results of using the original Breiman and Ihaka code. The *best reduced dimension* was obtained by evaluating models of all dimensions on the test data. The FDA/MARS procedure with degree = 2 and dimension 6 outperforms all of the others. In practice, test data might not be available to us, and we would need to use cross-validation or a similar technique to select this parameter.

Plotting the original data reveals that there is a strong person effect, reflected in both mean and variances irrespective of the vowel spoken. We standardized the data for each person so that each variable had mean zero and variance unity and redid the analysis. The lower part of Table 5 reports these improved results. Figure 4 shows some chosen projections of the 10-dimensional space of canonical variates found by the FDA/BRUTO procedure. Some classes stand out more distinctly than the corresponding plots for LDA shown in Figure 1.

4. MULTIRESPONSE ADAPTIVE REGRESSION PROCEDURES

In this section we define the class of multiresponse regression procedures that can be used and interpreted in the FDA proposal of Section 2 and give a number of examples. Following that, we provide details on the two particular choices that we focus on in this article—the MARS and BRUTO procedures.

Suppose that we have response variables Y_1, Y_2, \dots, Y_K and a vector of predictors \mathbf{X} . The class of regression models that we consider are of the form

$$\begin{aligned}\eta_1(\mathbf{X}) &= \beta_{01} + \sum_{m=1}^M \beta_{m1} h_m(\mathbf{X}), \\ \eta_2(\mathbf{X}) &= \beta_{02} + \sum_{m=1}^M \beta_{m2} h_m(\mathbf{X}), \\ &\vdots \\ \eta_K(\mathbf{X}) &= \beta_{0K} + \sum_{m=1}^M \beta_{mK} h_m(\mathbf{X}).\end{aligned}\quad (12)$$

The models for the different responses share a common set of basis functions but are allowed different coefficients β_{mk} .

Table 5. Results for Vowel Recognition Data

Technique	Error rates	
	Training	Test
1. LDA	.32	.56
Softmax	.48	.67
2. QDA	.01	.53
3. CART	.05	.56
4. CART (linear combination splits)	.05	.54
5. Single-layer perceptron		.67
6. Multilayer perceptron (88 hidden units)		.49
7. Gaussian node network (528 hidden units)		.45
8. Nearest-neighbor		.44
9. FDA/BRUTO	.06	.44
Softmax	.11	.50
10. FDA/MARS (degree = 1)	.09	.45
Best reduced-dimension (=2)	.18	.42
Softmax	.14	.48
11. FDA/MARS (degree = 2)	.02	.42
Best reduced-dimension (=6)	.13	.39
Softmax	.10	.50
12. Breiman and Ihaka	.16	.47
<i>Using standardized features</i>		
13. LDA	.25	.36
Softmax	.41	.60
14. QDA	.02	.53
15. CART	.04	.68
16. Nearest-neighbor		.46
17. FDA/BRUTO	.10	.29
18. FDA/MARS (degree = 1)	.11	.30
Softmax	.09	.47

These coefficients may be found by least squares, penalized least squares, or some other method. Finally, the basis functions $h_m(\mathbf{X})$ will often be chosen adaptively; to obtain the linearity required in the FDA procedure, the basis functions are treated as fixed after they are selected.

Besides the MARS and BRUTO procedures, this class includes most well-known multiresponse regression methods. Here are some examples:

- Multiresponse projection pursuit regression (SMART; Friedman 1984). In this model the K response variables share a common set of ridge functions, $h_m(\mathbf{X}) = h_m(\alpha_m^T \mathbf{X})$.
- Regression trees (CART; Breiman et al. 1984). A fitted regression tree can be represented by (12), with each h_m an indicator function of a rectangular region of predictor space.
- Neural networks. Consider a multilayered neural network with output functions linear in each of M hidden units in the top layer. The implicit model in such a network is of the form (12), where $h_m(\mathbf{X})$ is the output of the m th unit.
- Hinge functions. Breiman (1991a) described a multi-response version of his hinge function method that uses a common set of hinge functions for each response.

Operationally, any linear smoother $S(\lambda)$ could be used in the FDA procedure, not just penalized expansions onto basis functions as described here. But the latter allows a crisp interpretation as a PDA, whereas the former does not. In the remainder of this section we give details of MARS and BRUTO.

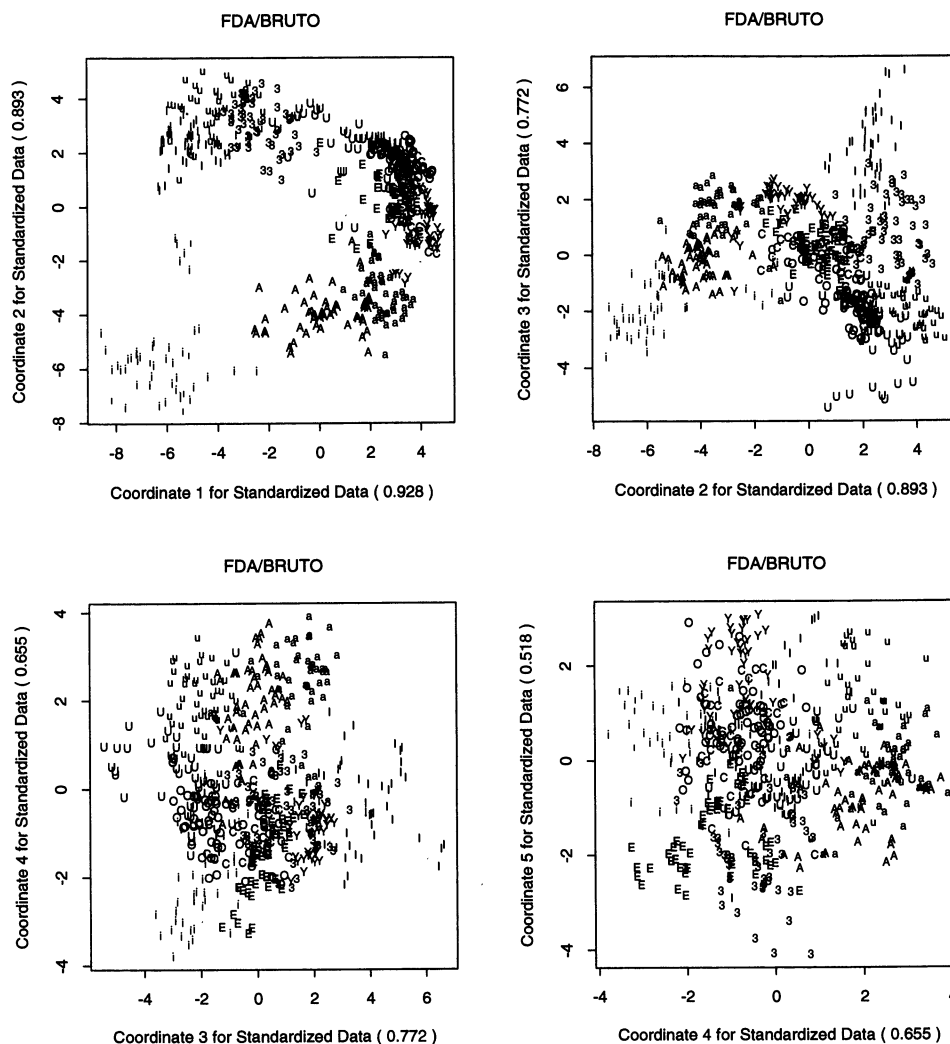


Figure 4. Some Selected Two-Dimensional Projections of the Canonical Variates Derived From the FDA/BRUTO Model Fit to the Vowel Data. The eigenvalues for each coordinate are in parentheses and indicate the strength of that variable.

4.1 MARS

The MARS proposal of Friedman (1991) is a procedure for adaptive nonparametric regression. It approximates the regression function by a model of the form

$$f(\mathbf{x}) = \beta_0 + \sum_{m=1}^M \beta_m \prod_{k=1}^{L_m} h_{lm}(x_{v(l,m)}). \quad (13)$$

Here x_1, x_2, \dots, x_p are predictors and $v(l, m)$ is the index of the predictor used in the l th term of the m th product. The basis functions h_{lm} are defined in pairs:

$$h_{lm}(x) = [x - t_{lm}]_+$$

and

$$h_{l,m+1}(x) = [t_{lm} - x]_+ \quad (14)$$

for m an odd integer, where the knot value t_{lm} is one of the unique values of $x_{v(l,m)}$. The model (13) is a sum of products of piecewise linear basis functions of the form (14). Denoting a typical model term by $H_j(\mathbf{x})$, the model is built up in a forward stepwise manner as follows:

1. We start with the constant basis function in the model.
2. At each state we consider adding to the model two terms defined by the products $H_j(\mathbf{x})h_{lm}(\mathbf{x})$ and $H_j(\mathbf{x})h_{l,m+1}(\mathbf{x})$ (m odd). If $H_j(\mathbf{x})h_{lm}(\mathbf{x})$ and $H_j(\mathbf{x})h_{l,m+1}(\mathbf{x})$ cause the greatest decrease in residual sum of squares, then they are added to the model. The forward stepwise process is stopped when some maximum model size is reached.
3. A backward "pruning" procedure is applied to the model, removing the least important terms one at a time. The best-fitting model in the stepwise sequence is chosen, with the fit measured by the generalized cross-validation (GCV) criterion.

The inclusion of each basis function and its conjugate pair in step 2 may cause redundancy in the basis but allows both members of the pair to be used in subsequent forward steps. The *degree* is a parameter of the procedure and limits the number of different terms in the product terms H_j and hence the order of interactions allowed. First-degree models are additive, second-degree models allow pairwise interactions, and so on. Full details of MARS may be found in Friedman (1991), who demonstrated the ability of MARS to capture

structure in high-dimensional data. The effectiveness of MARS stems mainly from:

- the piecewise linear basis functions, which permit fast updating of the least-squares fit as the knot position is changed, and
- the forward stepwise product strategy, which builds up the surface in small pieces. It takes advantage of any low-order structure that may exist and spends its degrees of freedom in a parsimonious manner.

We have generalized MARS for incorporation into the adaptive discriminant procedure. We coded an implementation of MARS from scratch, allowing for multiple response variables. In particular, for K response variables, K simultaneous models are estimated, having the same set of basis functions but different coefficients. The notation in (13) easily accommodates this, if each of the coefficients β_m is a K -vector. The models are built and pruned exactly as in Friedman's MARS, with the only difference being that the residual sum of squares and GCV criterion involve sums over the K response variables.

This generalization did not significantly complicate the implementation. In particular, Friedman's updating formulas still apply, allowing a rapid search for the knot locations.

When MARS is used in the FDA procedure, the relevant scoring criterion is

$$ASR = \frac{1}{N} \sum_{k=1}^K \sum_{i=1}^N \left(\theta_k(g_i) - \sum_{m=1}^M \beta_{mk} \prod_{l=1}^{K_m} h_{lm}(x_{v(l,m)}) \right)^2.$$

4.2 BRUTO: ADAPTIVE ADDITIVE MODELING

The model considered here is additive and thus more restrictive than the MARS model. On the other hand, it deals with the predictors in a smoother fashion and can be effective in situations with a large number of predictors.

We introduced the additive model in Section 2.6; for simplicity, here we will work with a single-response version of (6),

$$\sum_{i=1}^N \left(y_i - \sum_{j=1}^p \beta_j^T h_j(x_{ij}) \right)^2 + \sum_{j=1}^p \lambda_j \beta_j^T \Omega_j \beta_j. \quad (15)$$

Fortunately, (15) can be simplified, and the solution for fixed λ_j can be characterized by a simple set of estimating equations,

$$\mathbf{f}_j = \mathbf{S}_j(\lambda_j) \left\{ y_i - \sum_{k \neq j} f_k(x_{ik}) \right\}. \quad (16)$$

Here the \mathbf{S}_j denote smoothing spline operators for smoothing against the j th predictor. These equations lend themselves naturally to an iterative Gauss–Seidel algorithm, also known as backfitting (see Buja, Hastie, and Tibshirani 1989 for further details).

We still need to select which terms to include in the model, choose the values of λ for those terms selected, and modify the entire procedure to accommodate multiple responses. For a given set of p variables, we could optimize a criterion such as GCV over all p smoothing parameters λ_j in the p -

term additive model. Gu and Wahba (1988) described an efficient algorithm for doing this; however, their algorithm requires $O(N^3)$ computations, and so a cheaper approximation is attractive. The $O(N)$ algorithm outlined here was described by Hastie (1989); named BRUTO, it combines both backfitting and smoothing parameter selection.

The GCV for an additive model is defined as

$$GCV(\lambda) = \frac{\sum_{i=1}^N \{y_i - \sum_j \hat{f}_{j,\lambda_j}(x_{ij})\}^2}{N\{1 - \text{tr } \mathbf{S}(\lambda)/N\}^2}. \quad (17)$$

$\mathbf{S}(\lambda)$ is the *additive* operator for the given values of the smoothing parameters, and the terms \hat{f}_{j,λ_j} denote the fitted functions corresponding to these parameters. Computations of the $\text{tr}(\mathbf{S}(\lambda))$ in the denominator requires $O(N^3)$ operations. BRUTO attempts to minimize the GCV -like statistic

$$GCV^B(c, \lambda) = \frac{\sum_{i=1}^N \{y_i - \sum_j \hat{f}_{j,\lambda_j}(x_{ij})\}^2}{N(1 - [1 + c \cdot \sum_j df_j]/N)^2}, \quad (18)$$

where $df_j = \text{tr } \mathbf{S}_j(\lambda_j) - 1$ measures of the approximate degrees of freedom used in the j th smooth term and c is the cost per degree of freedom as described in Section 2.8. So $GCV^B(1, \lambda)$ approximates $\text{tr } \mathbf{S}(\lambda)$ by $1 + \sum_j \{\text{tr } \mathbf{S}_j(\lambda_j) - 1\}$, which requires only $O(N)$ computations.

We can write the criterion (18) in a form that focuses on the j th term:

$$GCV^B(1, \lambda) = \frac{\sum_{i=1}^N \{r_i^{(j)} - f_{j,\lambda_j}(x_{ij})\}^2}{N(1 - [df^{(j)} + \text{tr } \mathbf{S}_j(\lambda_j)/N])^2}, \quad (19)$$

where $r^{(j)} = y_i - \sum_{k \neq j} \hat{f}_k(x_{ij})$ is the partial residual in (16), fixing all the other terms in the model, and likewise $df^{(j)}$ are the degrees of freedom for the other terms, an *offset* in the GCV criterion. This suggests that we can minimize GCV^B iteratively, by focusing on one parameter at a time. Each step simply requires a univariate GCV optimizer, slightly modified to include the degrees of freedom for other terms in the model. Including the cost c adds no additional computational burden.

An important modification to the smoothers allows the null fit (mean) or the linear fit to be candidates. Fitting a term by a constant effectively *removes* it from the model, which already includes a constant. In this way the BRUTO algorithm includes variable selection in its model-selection procedure.

A final modification accommodates multiple responses. As shown in (6), multiple responses are accommodated in an efficient manner, using the same basis functions and penalty matrices. Instead of a coefficient vector, we have a coefficient matrix for each term. Similarly, the backfitting procedure generalizes, and with an omnibus GCV criterion, the entire procedure generalizes.

5. LINEAR DISCRIMINANT ANALYSIS VERSUS MULTIPLE LINEAR REGRESSION

Given a multiresponse regression procedure, there is a simpler way to perform classification, as mentioned in Section 2.2. The regression procedure is applied to an indicator response matrix \mathbf{Y} representing the classes, and a new observation is assigned to the class with largest fitted value. We

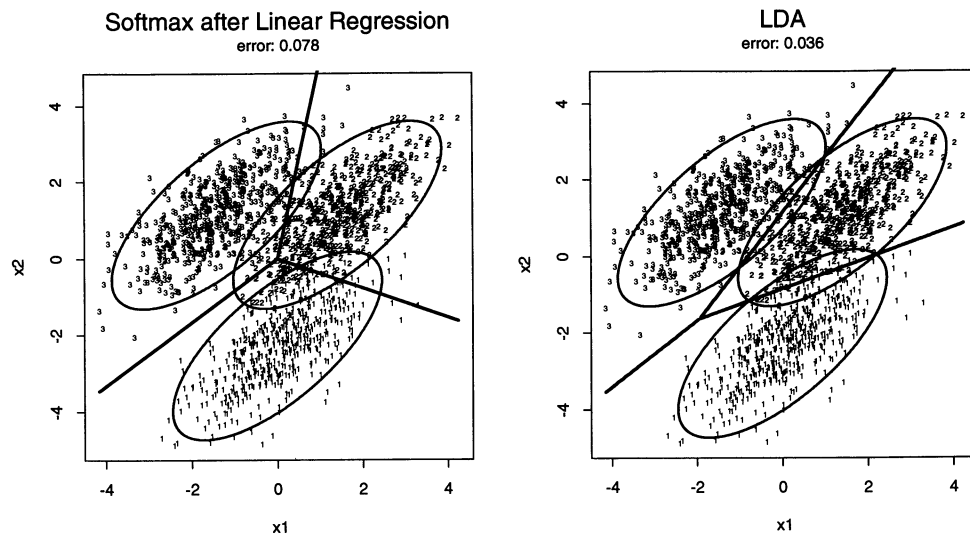


Figure 5. Softmax Uses the Wrong Metric. The data consists of 500 samples each from three bivariate Gaussian distributions, differing in location but with the same covariance structure. The location of their centroids distorts S_T relative to S_W and results in a biased decision boundary for softmax on the left but not for LDA on the right.

will refer to this procedure as *softmax*, as is done in the neural network literature. Notice that steps 1 and 2 of the LDA algorithm in Section 2 compute \hat{Y} ; so if we stopped there, the classification rule would be the same as that obtained from the *softmax* rule. Instead, we apply steps 3 and 4, then assign to the closest centroid in the η space. These remaining steps may be regarded as a postprocessor for converting the regression procedure into a discriminant analysis, and they have the following advantages:

1. Accuracy. We demonstrate that for three or more classes, *softmax* can produce ridiculous decision boundaries. It seems the fault lies in the fact that it uses the wrong metric in computing distances, as well as the wrong threshold.
2. Dimension reduction. Training data may produce some dimensions with marginal and most likely spurious class separations. Discrimination on test data can benefit from weeding out these weak dimensions (see the vowel data test results in Table 5.) Furthermore, the dimension reduction provides an important graphical tool for looking at the classification procedure when the number of classes is large.
3. Normal heuristics. The discriminant classifier is Bayes optimal for normally distributed data (assuming that class centers and shared class covariances are known). In practice the nonparametric canonical variates are linear combinations of many basis functions, and sums of random variables suggest normality.

On the other hand, Bayesian heuristics drive the *softmax* procedure. If the fitted values from the regression approximate the posterior probabilities well, then *softmax* will approximate the optimal Bayes procedure. Similar arguments can be used asymptotically for a consistent nonparametric regression model. Thus it seems that the success of FDA rests on the fact that the posterior probabilities are estimated with error by the regression method. We are currently studying this phenomenon in more detail.

In the preceding data examples, discriminant analysis almost always outperformed *softmax*, sometimes by a large

margin. In the remainder of this section, we try to gain some insight into the relationship between discriminant analysis and *softmax* after multiple regression. First, we view *softmax* as a minimum distance method, then we turn things around and view discriminant analysis as a regression method.

In the remainder of this section, by \mathbf{X} we mean either the original predictor matrix or else some adaptively chosen basis expansion of it. Without loss of generality, we assume that the columns of \mathbf{X} have mean zero. As before, \mathbf{m}_j is the mean vector of class j and S_T , S_{Bet} , and $S_W = S_T - S_{\text{Bet}}$ are the total, between-class, and within-class sample covariance matrices. The prior probability of class j is denoted by π_j .

5.1 Multiple Regression as a Distance Method

When the class priors are equal, Fisher's discriminant rule assigns to the class closest in Mahalanobis distance,

$$\delta_{\text{LDA}}(\mathbf{x}, j) = (\mathbf{x} - \mathbf{m}_j)^T S_W^{-1} (\mathbf{x} - \mathbf{m}_j). \quad (20)$$

We show that the corresponding *softmax* rule uses the distance

$$\delta_{\text{Soft}}(\mathbf{x}, j) = (\mathbf{x} - \mathbf{m}_j)^T S_T^{-1} (\mathbf{x} - \mathbf{m}_j) - \mathbf{m}_j^T S_T^{-1} \mathbf{m}_j. \quad (21)$$

They differ in the distance metric, and *softmax* has an additional threshold term.

Softmax uses the total covariance matrix to define the metric in computing distances, whereas LDA uses the within-covariance matrix. Because $S_T = S_W + S_{\text{Bet}}$, the configuration of the centroids can distort this metric. Figure 5 illustrates this phenomenon with an artificial example; we see the biased decision boundaries compared to those produced by LDA. This distortion can only occur when there are three or more groups. For two groups, it can be shown that the metric is the same in both cases and that only the threshold is different, by an amount $O((\pi_1 - 1/2)^3)$ (Ripley, personal communication; see also PDA.)

Softmax has an additional threshold term that also depends on the positions of the class centroids. If the classes

are all symmetrically placed about the origin, then this term plays no role. A worst-case situation for three classes is illustrated in Figure 6, where the class centroids are collinear. The middle class has centroid zero, while the outer classes do not, and each class benefits from the threshold correction in (21). (Think of the three dummy-variable regressions for more intuition in this case.)

An opposing argument in favor of *softmax* suggests that these phenomena are a result of the constraints imposed by the linear regression model. If the linear regression procedure is consistent for the underlying conditional expectations $p_j(\mathbf{x}) = E(Y_j|\mathbf{x})$, then *softmax* is the Bayes-optimal classifier. But this is unlikely to be the case unless we have sufficient data and thus sufficient basis functions.

We now derive (21). *Softmax* assigns to the class with the largest fitted value

$$R_{\text{Soft}}(j, \mathbf{x}) = \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_j + \frac{N_j}{N}. \quad (22)$$

This corresponds to minimizing the distance function,

$$\delta_{\text{Soft}}(x, j) = \frac{N_j}{2N} \{ (\mathbf{x} - \mathbf{m}_j)^T \mathbf{S}_T^{-1} (\mathbf{x} - \mathbf{m}_j) - \mathbf{x}^T \mathbf{S}_T^{-1} \mathbf{x} - \mathbf{m}_j^T \mathbf{S}_T^{-1} \mathbf{m}_j - 2 \}. \quad (23)$$

When all the N_j are equal, (23) is equivalent to (21).

5.2 Antipenalization

Here we view LDA as a regression method. One form of Fisher's linear discriminant rule is to assign to the class with the largest value of

$$R_{\text{LDA}}(j, \mathbf{x}) = \mathbf{x}^T \mathbf{S}_W^{-1} \mathbf{m}_j - \frac{1}{2} \mathbf{m}_j^T \mathbf{S}_W^{-1} \mathbf{m}_j + \log \pi_j. \quad (24)$$

The corresponding multiple regression rule (*softmax*) is to assign to the class with the largest fitted value

$$\begin{aligned} R_{\text{MR}}(j, \mathbf{x}) &= \mathbf{x}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}_j + N_j / N \\ &= N_j / N (\mathbf{x}^T \mathbf{S}_T^{-1} \mathbf{m}_j + 1). \end{aligned} \quad (25)$$

Again, to simplify the comparison, we consider the equal N_j situation and focus on the *class coefficient vectors*. These are up to a constant,

$$b_j^{\text{LDA}} = \mathbf{S}_W^{-1} \mathbf{m}_j \quad (\text{LDA})$$

$$b_j^{\text{Soft}} = \mathbf{S}_T^{-1} \mathbf{m}_j \quad (\text{softmax}).$$

It is not hard to show that these solve the minimization problems:

$$\frac{1}{N} \|\mathbf{y}_j - \mathbf{X} \mathbf{b}_j\|^2 - \mathbf{b}_j^T \mathbf{S}_{\text{Bct}} \mathbf{b}_j \quad (\text{LDA}) \quad (26)$$

$$\frac{1}{N} \|\mathbf{y}_j - \mathbf{X} \mathbf{b}_j\|^2 \quad (\text{softmax}). \quad (27)$$

Because the penalty term in (26) is subtracted rather than added, we call this *antipenalization*. These equations offer some further insight into the comparison of LDA and *softmax*. In *softmax*, the direction b_j is chosen to separate out class j from the rest. In contrast, the extra term in the antipenalization criterion (26) rewards directions b_j that also separate out all of the J classes.

This comparison has some caveats:

- When the N_j are not equal, this interpretation gets a bit murky.
- We have ignored the constant terms, which could be included in both criteria—without penalty in (26). The *antipenalized* criterion produces a constant different from that in (24)—in fact, different essentially by the amount $\frac{1}{2} \mathbf{m}_j^T \mathbf{S}_W^{-1} \mathbf{m}_j$. A similar term showed up as crucial in the second example in the previous section.

5.3 Other Postprocessors

Some other “postprocessors” have been proposed; for example, as nearest-neighbor classification in the J -dimensional space. (For this and other approaches in the context of neural networks, see Denker and LeCun 1990.) Another approach is to use the basis functions from step 2 in a multinomial

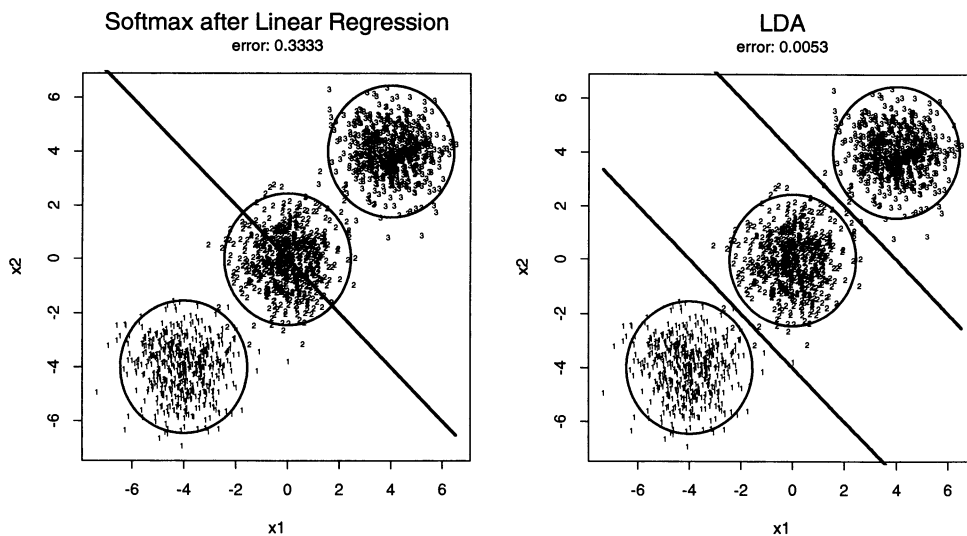


Figure 6. Softmax Uses the Wrong Threshold. The data consist of 500 samples each from three spherical bivariate Gaussian distributions, whose centroids line up along a line. The center class is completely masked by the outside two when softmax is used, but LDA has no such problem.

regression. Results in standard (linear) logistic and multinomial regression (Bull and Donner 1987; Efron 1975) suggest that this might be more robust to nonnormality and might improve classification performance. But because the normality refers here to the space of adaptively chosen basis functions (not the original variables), it is not clear how the methods would compare.

6. FURTHER DETAILS AND DISCUSSION

6.1 Iterative FDA

We can see at least two situations where iterating the FDA algorithm could lead to improvements:

- For $K < J - 1$ -dimensional models, iteration can lead to better choices of Θ_0 .
- Ideally, we would like the criterion for choosing λ to reflect misclassification errors; in iterative procedure emerges from such considerations.

6.1.1 Low-rank Models. In Section 2.7 we saw that the *ASR* and hence the *GCV* criterion are invariant under different initial scores Θ_0 of rank $J - 1$, but not necessarily so for lower rank. Sometimes lower-rank models, being more parsimonious and hence less variable, can lead to better classification results (see the vowel data results, for example). One approach is to use a full-rank Θ_0 and then extract a K -dimensional subspace (which is what we did for the vowel example.) It is conceivable, however, that the search for a suitable λ might benefit if a more focused *ASR* criterion were used—one that depends only on the scorings in question. Because these are not known in advance, an iteration is suggested. We are currently exploring such iterative approaches, especially for problems with large J .

6.1.2 Classification Criterion. In Section 2.8 we saw that our criterion for selecting λ is most conveniently a regression criterion, but that our real focus is on classification errors and group separation. The following identity suggests a way of bringing the two goals closer together. The (penalized) Mahalanobis distance from an observation x to the j th group is given by

$$\delta_\lambda(\mathbf{x}, j) = \sum_{k=1}^J \frac{(\eta_k(\mathbf{x}) - \theta_{jk})^2}{1 - \alpha_k^2} - 1/p_j + C(h(\mathbf{x})), \quad (28)$$

where θ^j is the j th row of Θ . This identity is due to Breiman and Ihaka (1984); in PDA we derive the more general version that covers penalized distances.

Mahalanobis distance is one step away from classification, because we assign classes to the closest group. This suggests using a weighted sum-of-squares or *discriminant distance* criterion,

$$DD(\lambda) = \sum_{i=1}^N \sum_{k=1}^{J-1} \frac{(\eta_k(\mathbf{x}_i; \lambda) - \theta_k(g_i; \lambda))^2}{1 - \alpha_k^2(\lambda)}, \quad (29)$$

as the basis for finding λ . This weighted sum-of-squares criterion depends on Θ (and α) and suggests an obvious alternating algorithm.

We can only report empirical evaluations of this approach. On all the examples tried, the relative class distance was

decreased on average over the training data. This resulted in a modest decrease in the average misclassification errors over the training data, although not always for the test data.

6.2 Priors and Posteriors

So far we have ignored the role of class prior probabilities. Priors can play a role in two places in our procedure:

1. The Bayes classification procedure makes use of the priors in minimum cost assignment.
2. Priors can be used in estimation to weight the contributions of observations.

With a 0/1 loss function, the Bayes-optimal procedure is to classify an observation x to the class with highest posterior probability. For the Gaussian model, the posterior for the j th class is given by

$$p_j(x) \propto \pi_j \exp[-(\mathbf{x} - \mathbf{m}_j)^T \Sigma_w^{-1} (\mathbf{x} - \mathbf{m}_j)/2] \\ \propto \exp[\{-\delta(\mathbf{x}, j) - 2 \log \pi_j\}/2], \quad (30)$$

where π_j is the prior for class j . So classifying to the class with the minimum discriminant score $\delta(\mathbf{x}, j)$ assumes that the priors are equal; if the priors are known and unequal, then posteriors can be computed as in (30) or else the distances can be corrected by subtracting $2 \log(\pi_j)$. If the training sample is random, then it is natural to estimate the π_j by the empirical priors N_j/N .

Our estimation procedure implicitly weights each class by the sample prior, because each class is represented by that many observations. This is made more explicit in our eigen-decomposition, when we use D_p as the normalization metric for the scores. If the sample is stratified and the true priors are known, then it might be more appropriate to weight the observations differently to reflect this knowledge, using weights $N\pi_j/N_j$ for observations in the j th class and replacing D_p by D_π .

6.3 Discussion

Breiman and Ihaka (1984) suggested a procedure like the one described in this article, fitting an additive model to the columns of Θ_0^* for a specific choice of Θ_0 . In their public domain software (available from the *statlib* archive at Carnegie-Mellon University), they fitted a separate additive model to each set of scores, using an adaptive procedure based on "supersmoother" (Friedman and Stuetzle 1982). In addition, they applied the adaptive fitting procedure in step 4 as well as step 2. We have chosen to fit a common model to the set of $J - 1$ scores simultaneously; for our procedure, refitting in step 4 is not necessary. Our reasons for doing this are as follows:

1. Our procedure has a clear interpretation as a linear (or penalized) LDA on a set of adaptively chosen functions of the predictors; this interpretation is not valid for their procedure.
2. Fitting one model rather than K models requires roughly $1/K$ times as much computer time. Although speed considerations should not be of primary concern, it is a factor, especially when full cross-validation is used.

3. Fitting K separate models uses up too many degrees of freedom, too quickly.

On the other hand, some problems might benefit from the extra flexibility from having a separate model per dimension. Our experience with the separate approach is that it is generally (slightly) inferior in performance (e.g., Table 5), and hence does not warrant the extra complexity.

Software. We have written a general procedure in the S language (Becker and Chambers 1988) called `fda` for FDA. A `method` argument allows the user to specify the multi-response regression method to be used; the default is linear regression and thus Fisher's LDA. The regression method must be provided in a separate S function. For example, `method = mars` or `method = bruto` use the MARS and BRUTO methods respectively. The software will be publicly available in the statistics archive at Carnegie-Mellon University (`statlib@lib.stat.cmu.edu`)

APPENDIX: EQUIVALENCE OF PENALIZED DISCRIMINANT ANALYSIS AND PENALIZED OPTIMAL SCORING

Here we present a terse proof of the equivalence; a more thorough treatment is given in PDA. Our proof includes the equivalence of LDA and optimal scoring as a special case, simply by setting the penalty Ω to zero. Let Σ be the sample covariance matrix of $(Y : H)$, with Σ_{ij} , $i, j \in \{1, 2\}$ the obvious partitions. As before, we assume that the columns of H are centered and that H represents a basis expansion of the original variables (or in the unpenalized case X itself.) The following connections exist with our previous notation:

- $\Sigma_{11} = D_p$, the diagonal matrix of class proportions
- $\Sigma_{11}^{-1} \Sigma_{12} = M$, the matrix of class centroids
- $\Sigma_{22} = \Sigma_T$
- $\Sigma_{\text{Bet}} = \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12}$.

Consider the *generalized, penalized, singular-value decomposition* of Σ_{12} :

$$\Sigma_{11}^{-1} \Sigma_{12} (\Sigma_{22} + \Omega/N)^{-1} = \Theta_* D_\alpha B_*^T, \quad (\text{A.1})$$

with normalizations

$$\Theta_*^T \Sigma_{11} \Theta_* = I = B_*^T (\Sigma_{22} + \Omega/N) B_*. \quad (\text{A.2})$$

We now show that, modulo constants, these same left and right singular matrices diagonalize both the penalized optimal scoring criterion and the PDA criterion and hence are optimal for both.

A.1 Penalized Optimal Scoring

We start with the partially optimized criterion (7),

$$ASR_p(\Theta) = \text{tr}\{\Theta^T Y^T (I - S) Y \Theta\} / N.$$

Substituting the appropriate components of Σ , we see that

$$\begin{aligned} ASR_p(\Theta_*) &= \text{tr}\{\Theta_*^T (\Sigma_{11} - \Sigma_{12} (\Sigma_{22} + \Omega/N)^{-1} \Sigma_{21}) \Theta_*\} \\ &= \text{tr}(I - D_\alpha^2). \end{aligned}$$

The penalized regression coefficients are given by

$$B = (\Sigma_{22} + \Omega/N)^{-1} \Sigma_{21} \Theta_* = B_* D_\alpha.$$

A.2 Penalized Discriminant Analysis

One statement of the PDA problem is to maximize $\text{tr}\{U^T \Sigma_{\text{Bet}} U\}$ subject to the penalized constraint $U^T (\Sigma_W + \Omega/N) U = I$. Now

$$B_*^T \Sigma_{\text{Bet}} B_* = B_*^T \Sigma_{21} \Sigma_{11}^{-1} \Sigma_{12} B_* = D_\alpha^2$$

and

$$B_* (\Sigma_W + \Omega/N) B_* = B_*^T (\Sigma_{22} + \Omega/N - \Sigma_{\text{Bet}}) B_* = I - D_\alpha^2.$$

Thus $U = B_* (I - D_\alpha^2)^{-1/2}$ satisfies the constraint and diagonalizes the criterion, and hence is the optimal discriminant coefficient matrix.

Connecting these two results, we see that $U = BD$, where

$$D^{-1} = D_\alpha (I - D_\alpha^2)^{1/2}.$$

[Received February 1993. Revised December 1993.]

REFERENCES

- Barron, A. R., and Barron, R. L. (1988), "Statistical Learning Networks: A Unifying View," in *Computer Science and Statistics: Proceedings of the 21st Interface*.
- Becker, R., Chambers, J., and Wilks, A. (1988), *The New S Language*, Belmont, CA: Wadsworth.
- Breiman, L. (1991a), "Hinging Hyperplanes for Regression, Classification and Function Approximation," Technical Report 324, University of California, Berkeley, Dept. of Statistics.
- (1991b), "The π -Method for Estimating Multivariate Functions From Noisy Data," *Technometrics*, 33(2), 125–160.
- Breiman, L., and Friedman, J. (1985), "Estimating Optimal Transformation for Multiple Regression and Correlation" (with discussion), *Journal of the American Statistical Association*, 80, 580–619.
- Breiman, L., Friedman, J. H., Olshen, R., and Stone, C. J. (1984), *Classification and Regression Trees*, Belmont, CA: Wadsworth.
- Breiman, L., and Ihaka, R. (1984), "Nonlinear Discriminant Analysis via Scaling and ACE," technical report, University of California, Berkeley, Dept. of Statistics.
- Bridle, J. S. (1989), "Training Stochastic Model Recognition Algorithms as Networks can Lead to Maximum Mutual Information Estimation of Parameters," *Advances in Neural Information Processing Systems*, Touretzky, D. ed., 2.
- Buja, A., Hastie, T., and Tibshirani, R. (1989), "Linear Smoothers and Additive Models" (with discussion), *The Annals of Statistics*, 17, 453–555.
- Bull, S., and Donner, A. (1987), "The Efficiency of Multinomial Logistic Regression Compared With Multigroup Normal Discriminant Analysis," *Journal of the American Statistical Association*, 82, 1118–1122.
- Denker, J. S., and Le Cun, Y. (1991), "Transforming Neural-Net Output Levels to Probability Distributions," *Advances in Neural Information Processing Systems*, Vol. 3, eds. R. Lippmann, J. Moody, and D. Touretzky, Denver: Morgan Kaufman.
- Efron, B. (1975), "The Efficiency of Logistic Regression Compared to Normal Discriminant Analysis," *Journal of the American Statistical Association*, 70, 892–898.
- Friedman, J. (1987), "Exploratory Projection Pursuit," *Journal of the American Statistical Association*, 82, 249–266.
- (1991), "Multivariate Adaptive Regression Splines" (with discussion), *The Annals of Statistics*, 19, 1–141.
- Friedman, J., and Stuetzle, W. (1981), "Projection Pursuit Regression," *Journal of the American Statistical Association*, 76, 817–823.
- (1982), "Supersmoother," technical report, Stanford University, Dept. of Statistics.
- Gifi, A. (1981), "Nonlinear Multivariate Analysis," unpublished manuscript.
- (1990), *Nonlinear Multivariate Analysis*, Wiley, Chichester.
- Gnanadesikan, R., and Kettenring, J. (eds.) (1989), "Discriminant Analysis and Clustering," *Statistical Science*, 4, 34–69.
- Gorman, R., and Sejnowski, T. (1988), "Analysis of Hidden Units in a Layered Network Trained to Classify Sonar Targets," *Neural Networks*, 1, 75–89.
- Gu, C., and Wahba, G. (1984), "Minimizing GCV/GML Scores with Multiple Smoothing Parameters via Newton's Method," Technical Report 847, University of Wisconsin-Madison, Dept. of Statistics.
- Hastie, T. (1989), Discussion of "Flexible parsimonious smoothing and additive modelling" by Friedman and Silverman, *Technometrics*, 31, 3–39.
- Hastie, T., Buja, A., and Tibshirani, R. (1994), "Penalized discriminant analysis," submitted to *The Annals of Statistics*.
- Hastie, T., and Tibshirani, R. (1985), Discussion of "Projection Pursuit" by Huber, *The Annals of Statistics*, 13, 502–507.
- (1990), *Generalized Additive Models*, London: Chapman and Hall.

- Hinton, G. (1989), "Connectionist Learning Procedures," *Artificial Intelligence*, 40, 185-234.
- Lippman, R. (1989), "Pattern Classification Using Neural Networks," *IEEE Communications Magazine*, 11, 47-64.
- Loh, W., and Vanichsetakul, N. (1988), "Tree-Structured Classification via Generalized Discriminant Analysis," *Journal of the American Statistical Association*, 83, 715-728.
- Mardia, K., Kent, J., and Bibby, J. (1979), *Multivariate Analysis*, New York: Academic Press.
- Owen, A. (1991), Discussion of "Multivariate Adaptive Regression Splines" by J. Friedman, *The Annals of Statistics*, 19, 102-112.
- Rabiner, L. R., and Schafer, R. W. (1978), *Digital Processing of Speech Signals*, Englewood Cliffs, NJ: Prentice-Hall.
- Ripley, B. (1994), "Neural Networks and Related Methods for Classification," (with discussion), Vol. 56, No. 3, 409-456.
- Ripley, B., and Hjovt, N. (in press), *Pattern Recognition and Neural Networks—A Statistical Approach*, Cambridge: Cambridge University Press.
- Robinson, A. (1989), "Dynamic Error Propagation Networks," Ph.D. thesis, Cambridge University, Dept. of Electrical Engineering.
- Wahba, G. (1990), *Spline Models for Observational Data*, Philadelphia: SIAM.