

1 International Journal of Pattern Recognition
and Artificial Intelligence
3 Vol. 21, No. 7 (2007) 1–24
© World Scientific Publishing Company



5 **FEATURE EXTRACTION FOR CLASSIFICATION USING**
STATISTICAL NETWORKS

7 ANIL KUMAR GHOSH
Centre for Mathematics and Its Applications, Mathematical Sciences Institute
9 *Australian National University, Canberra, ACT 0200, Australia*
anilkghosh@rediffmail.com

11 SMARAJIT BOSE
Theoretical Statistics and Mathematics Unit, Indian Statistical Institute
13 *203, Barrackpore Trunk Road, Kolkata 700108, India*
smarajit@isical.ac.in

15 In a classification problem, quite often the dimension of the measurement vector is large.
17 Some of these measurements may not be important for separating the classes. Removal
of these measurement variables not only reduces the computational cost but also leads
19 to better understanding of class separability. There are some methods in the existing
literature for reducing the dimensionality of a classification problem without losing much
21 of the separability information. These dimension reduction procedures usually work well
for linear classifiers. In the case where competing classes are not linearly separable, one
23 has to look for ideal “features” which could be some transformations of one or more
measurements. In this paper, we have tried to tackle both the problems of dimension
25 reduction and feature extraction by considering a projection pursuit regression model.
The single hidden layer perceptron model and some other popular models can be viewed
27 as special cases of this model. An iterative algorithm based on backfitting is proposed
to select the features dynamically and cross-validation is used to select the ideal number
29 of features. We carry out an extensive simulation study to show the effectiveness of this
fully automatic method.

31 *Keywords:* Artificial neural networks; backfitting; classification using splines; cross-
validation; feature selection; projection pursuit regression.

1. Introduction

33 In high-dimensional classification problems, it is often seen that only a smaller
number of features (some functions of original variables) contain most of the infor-
35 mation for separating the classes. It is therefore desirable to identify these features
so that it helps to reduce the dimensionality of the problem which in turn facilitates
37 the visualization of class separability in a lower dimensional space and at the same
time reduces the computational and storage costs substantially.

2 *A. K. Ghosh & S. Bose*

1 If we restrict ourselves only to linear classifier, results for feature extraction are
3 available in the literature (see e.g. Refs. 14 and 24). However, in practice, linear clas-
5 sifiers are often found to be inadequate. Therefore, extraction of important features
7 for nonlinear classification is an important statistical problem. For linear classifiers
9 the features are typically linear combinations of the original variables, which we will
11 consider as linear features in this article. For nonlinear classification, in addition,
13 we consider features which are nonlinear transformations either of original variables
15 or of linear features. We will consider these as nonlinear features.

17 In regression, projection pursuit regression model (see e.g. Refs. 16 and 26)
19 was developed to compensate for this inadequacy, and artificial neural network
21 models^{8,34,35} provide much more flexibility for classification. Both have their own
23 strengths and limitations. In this paper, we try to combine these two approaches
25 so that we can develop a method which is more flexible and statistically meaning-
ful as well. It is shown later that many of the existing models are special cases
of the resulting model. We propose an iterative feature selection (IFS) algorithm
for dynamic adjustment of features and cross-validation techniques are used for
selection of the final model appropriate to a specific problem.

The organization of the paper is as follows. The problem of classification is
briefly described and the available tools for solving these problems are reviewed in
Sec. 2. In Sec. 3, we discuss the motivation for the new method and compare the
corresponding model with the existing ones. The algorithm is described in Sec. 4.
Experimental results illustrating the new method is presented in Sec. 5 where we
use some simulated and benchmark data sets to compare the method with some
existing classification techniques. Finally, a summary and concluding remarks are
given in Sec. 6.

2. Classification Techniques

27 In classification problems, one aims to develop a decision rule $d(\mathbf{x}) : R^d \rightarrow$
29 $\{1, 2, \dots, J\}$ based on the available training sample $\{(\mathbf{x}_n, c_n) : \mathbf{x}_n \in R^p, c_n \in$
31 $\{1, 2, \dots, J\}; n = 1, 2, \dots, N\}$ for classifying multivariate observations \mathbf{x} into one
33 of J competing populations. Bayes rule (see e.g. Ref. 1) assigns an observation \mathbf{x} to
35 the class with the largest conditional probability $p(j|\mathbf{x})$. An accurate approxima-
37 tion of these unknown probabilities helps us to formulate a decision rule capable
39 of achieving nearly the optimal misclassification rate. These probabilities can be
41 estimated either parametrically or nonparametrically. In parametric approaches
(see e.g. Refs. 1, 19 and 31), the measurement vector \mathbf{x} is assumed to have a
known distribution with unknown parameters. For instance, in Fisher's (see Ref. 15)
linear or quadratic discriminant analysis (LDA and QDA), this known distribu-
tion is taken as multivariate normal with different parameters for different classes.
Performance of a parametric classifier depends heavily on the validity of paramet-
ric model assumptions, and when one or more of these conventional assumptions
are violated, parametric classifiers often fail to adequately approximate the class

1 boundaries. In such situations, one usually prefers nonparametric methods, which
 2 are more flexible and robust. Notably nonparametric methods like classification
 3 trees (see e.g. Ref. 5), flexible discriminant analysis (see e.g. Ref. 23), nearest neigh-
 4 bors (see e.g. Refs. 10 and 18), kernel discriminant analysis (see e.g. Refs. 13 and
 5 22), neural networks and support vector machines (see e.g. Refs. 11, 37 and 42) have
 6 been shown to outperform the parametric approaches in a wide variety of problems.

7 Many of these classification techniques including the neural network method
 8 regress indicator variables Y_1, \dots, Y_J (defined for J competing classes) on the mea-
 9 surement space to estimate $E(Y_j|\mathbf{x}) = p(j|\mathbf{x})$ ($j = 1, 2, \dots, J$), the posterior class
 10 probabilities. In particular, Bose³ followed this principle to develop Classification
 11 Using Splines (CUS), which is based on nonparametric additive regression. This
 12 method essentially considers smooth functions of the original variables as new fea-
 13 tures and considers a linear model of these features to estimate posterior class prob-
 14 abilities. The proposed method in this paper can be viewed as a more generalized
 15 version of CUS.

3. Features for Nonlinear Classification

17 In classification problems, when class boundaries are highly nonlinear, one popu-
 18 lar practice is to find suitable nonlinear features such that the competing classes
 19 are linearly separable in that feature space. Linear classification in that feature
 20 space essentially leads to a nonlinear separation in the space of the original mea-
 21 surement variables. Well-known methods like support vector machines and kernel
 22 fisher discriminants (see e.g. Ref. 32) adopt these techniques and use some suitable
 23 nonlinear kernel functions for nonlinear classification. Ghosh and Chaudhuri²¹ also
 24 used a similar strategy for depth based nonlinear classification. However, instead
 25 of estimating features from the data, they worked with some predefined features.
 26 In a recent work, Zhu⁴³ tried to extract important features from the data by max-
 27 imizing a likelihood ratio criterion. Neural network training algorithms also try to
 28 extract ideal features from the data automatically, but this popular method lacks
 29 meaningful statistical interpretations.

30 As we have mentioned in the previous section, CUS considers additive models
 31 (see e.g. Ref. 7) to estimate the posterior probabilities, and the estimates are of
 32 the form

$$33 \quad \hat{p}(j|\mathbf{x}) = \phi_{j0} + \sum_{i=1}^p \phi_{ji}(x_i), \quad j = 1, 2, \dots, J, \quad (1)$$

34 where ϕ_{j0} is a constant and ϕ_{ji} are smooth univariate functions. Therefore, CUS
 35 may not perform well when the class boundaries cannot be approximated well by
 36 additive models.

37 Method of successive projection (see Ref. 4) was proposed to compensate for
 this inadequacy. It uses the estimated posterior probabilities obtained by CUS as
 the new measurement variables and repeats the CUS procedure with them to get

4 *A. K. Ghosh & S. Bose*

1 the final probability estimates

$$\hat{p}(j|\mathbf{x}) = \phi_{j0}^{(2)} + \sum_{k=1}^{J-1} \phi_{jk}^{(2)} \left(\phi_{k0}^{(1)} + \sum_{i=1}^p \phi_{ki}^{(1)}(x_i) \right), \quad j = 1, 2, \dots, J, \quad (2)$$

3 where $\phi^{(1)}$ and $\phi^{(2)}$ are smooth univariate functions estimated in the first and
second iterations, respectively.

5 In the process, interactions are introduced indirectly in the model much in the
same manner as in the multilayer perceptron model. Successive projection has two
7 basic steps, one for variable transformation and the other for probability estimation.
This is similar to the Multilayer Perceptron (MLP) model with a single hidden layer
9 where posterior estimates are of the form:

$$\hat{p}(j|\mathbf{x}) = \rho^{(2)}(\theta_j^{(2)}) + \sum_k w_{jk}^{(2)} \rho^{(1)} \left(\theta_k^{(1)} + \sum_m w_{mk}^{(1)} x_m \right), \quad (3)$$

11 where $\theta^{(r)}$, $w^{(r)}$ and $\rho^{(r)}$ denote bias terms, weight functions and transformation
functions associated with the r th ($r = 1, 2$) layer. Several choices for the transfor-
13 mation function ρ is available in the literature. Because of its ability for uniform
approximation of any continuous function (see e.g. Ref. 25) on a compact support,
15 sigmoidal function $\sigma(t) = 1/(1 + e^{-t})$ is one of the most popular choices.

17 However, one disadvantage of successive projection is that the transformations
are estimated once and for all. It lacks the flexibility of neural networks where
features are modified iteratively. As a result, when the posterior probability esti-
19 mates obtained by CUS are not good, the improvement over CUS is often not that
significant.

21 The proposed iterative feature selection (IFS) algorithm, instead of using CUS,
in the first step, tries to extract the best possible linear features (linear combina-
23 tions) of the form $\alpha'_i \mathbf{x}$ ($i = 1, 2, \dots, p$) from the data, and then performs CUS
with them to estimate the posterior class probabilities. Backfitting is used for
25 dynamic adjustment of α_i and that makes the method very flexible. Like neu-
ral networks, IFS can adjust the linear combinations $\alpha'_i \mathbf{x}$ iteratively to improve its
27 performance. It can be viewed as an approach to estimate the posterior probabilities
 $p(j|\mathbf{x})$, ($j = 1, 2, \dots, J$) by projection pursuit regression models

$$\hat{p}(j|\mathbf{x}) = \psi_j(\mathbf{x}) = \phi_{j0} + \sum_{i=1}^p \phi_{ji}(\alpha'_i \mathbf{x}). \quad (4)$$

29 In projection pursuit, one has to find the ideal linear combinations and the
31 corresponding transformation functions as well. IFS tries to accomplish that task
automatically, and the resulting decision rule seems to provide a much better
33 improvement over CUS, much more than the method of successive projections.

35 The model used in IFS has some similarities with some of the existing classi-
fication methods. Obviously, the additive model used in CUS is a special case of
IFS when α_i s are unit vectors along the co-ordinate axes. In fact, if sigmoidal or

1 other known transformations are used as the smooth univariate functions ϕ_{ji} , this
 2 model turns out to be a perceptron model with one hidden layer. In the perceptron
 3 model, a fixed set of transformation functions (which could possibly be different) is
 4 used, and a given function is approximated to any level of accuracy by increasing
 5 the number of hidden nodes. This implies more number of linear features (linear
 6 combinations) and hence a transformation to a high-dimensional feature space from
 7 the original measurement space.

8 Since our aim is to reduce the dimensionality of the problem, in IFS, we restrict
 9 the number of features to the dimension of the original measurement vector and
 10 try to approximate the transfer functions by using additive splines. Therefore, not
 11 only the linear combinations are estimated from the data iteratively but so are
 12 the transfer functions. The fewer number of linear combinations may help to identify
 13 simpler structures in the data which has been illustrated with the simulated
 14 Example 1 in Sec. 5.1.

15 Use of splines as transformation functions makes the posterior estimates affine-
 16 invariant, and it also leads to the universal approximation property like the
 17 sigmoidal function. Given the set of linear features, using sufficiently many knots,
 18 splines can approximate any continuous function on that space to any desired level
 19 of accuracy (see Ref. 3). A detailed discussion is given in Sec. 4.5. The IFS algo-
 20 rithm actually evolved from a similar scheme applied for training perceptrons ear-
 21 lier which resulted in Backfitting Neural Networks (see Ref. 20). Mackay²⁹ has also
 22 proposed many modifications to the traditional backpropagation training algorithm
 23 including regularizations, model selections but he considered a different Bayesian
 24 approach.

25 4. Description of IFS Algorithm

IFS projects the posterior probabilities $p(j|\mathbf{x})$, ($j = 1, 2, \dots, J$) into the additive
 class of functions spanned by the cubic spline functions of $\boldsymbol{\alpha}'_i \mathbf{x}$ ($i = 1, 2, \dots, p$).
 Hence $\hat{p}(j|\mathbf{x})$ is of the form:

$$\hat{p}(j|\mathbf{x}) = \psi_j(\mathbf{x}) = \phi_{j0} + \sum_{i=1}^p \phi_{ji}(\boldsymbol{\alpha}'_i \mathbf{x}), \quad \text{where } \|\boldsymbol{\alpha}_i\| = 1 \quad \forall i = 1, 2, \dots, p. \quad (5)$$

26 The constants ϕ_{j0} , direction vectors $\boldsymbol{\alpha}_i$ and the transformation functions (cubic
 27 splines) ϕ_{ji} ($j = 1, 2, \dots, J$; $i = 1, 2, \dots, p$) are determined iteratively. Cubic splines
 28 are pieces of cubic polynomials joined together at the knots suitably placed on the
 29 range of a variable. They are sufficiently smooth, and have continuous derivatives up
 30 to the second order. Following the suggestion of Breiman,⁷ we also put restrictions
 31 on splines to make the extrapolated fit linear in the tails for reducing the effect of
 32 high variability at the end points. IFS starts with sufficiently many knots and then
 33 at the end of the iterative learning procedure, it reaches a smaller dimensionality
 34 by the method of backward deletion. The optimum dimensionality is selected by
 35 cross-validation. Guidelines for other related issues like selection of basis functions,

6 A. K. Ghosh & S. Bose

1 method of backward deletion are as given in Ref. 7 which we discuss briefly in
this section.

3 4.1. Selection of basis functions and parameters

For data with univariate $Z_n = \boldsymbol{\alpha}'\mathbf{x}_n$, ($n = 1, 2, \dots, N$) and K knots $t_1 < t_2$
5 $< \dots < t_K$ in the range [$Z_{(1)} = \min_n Z_n, Z_{(N)} = \max_n Z_n$], a cubic spline ϕ is a
cubic polynomial in each interval (t_k, t_{k+1}) , for $k = 1, 2, \dots, K - 1$, and it has con-
7 tinuous derivatives up to the second order. It is easy to observe that for a given set
of K knots, the space of cubic splines is $(K + 4)$ dimensional. Two sets of functions
9 commonly used as the bases for the class of cubic splines are the power basis and the
B-spline basis (see e.g. Ref. 12). In our algorithm, we use the former because dele-
11 tion of knots is computationally more feasible with this basis. The power basis for
univariate splines consists of the functions $1, Z, Z^2, Z^3, (Z - t_k)_+^3$, $k = 1, 2, \dots, K$,
13 where $Z_+ = Z$ if $Z > 0$ and 0 otherwise. To make the extrapolated spline fit
outside the range [$Z_{(1)}, Z_{(N)}$] linear, Breiman⁷ suggested to impose the conditions
15 $\phi''(Z_{(1)}) = \phi'''(Z_{(1)} -) = 0$ for the left tail, and $\phi''(Z_{(N)}) = \phi'''(Z_{(N)} +) = 0$ for the
right tail.

17 A convenient power basis for this space can be constructed by starting with the
functions $1, Z, (Z - t_k)_+^3$, $k = 1, 2, \dots, K$. Therefore, the restricted spline function
19 is of the form

$$\phi(Z) = \theta + \gamma Z + \sum_{k=1}^K \beta_k (Z - t_k)_+^3, \quad (6)$$

21 which automatically satisfies the condition for the left tail. To make $\phi(Z)$ linear
to the right of $Z_{(N)}$, Breiman⁷ suggested to take $t_1 = Z_{(1)}$ and an additional knot
23 t_{K+1} at $Z_{(N)}$ to cancel out the higher order terms. Condition for linearity at the
right end imposes the constraints

$$\sum_{k=1}^{K+1} \beta_k = 0 \quad \text{and} \quad \sum_{k=1}^{K+1} \beta_k (Z - t_k) = 0 \quad \text{for} \quad Z \geq Z_{(N)}. \quad (7)$$

25 For multivariate $\mathbf{Z} = (Z_1 = \boldsymbol{\alpha}'_1 \mathbf{x}, Z_2 = \boldsymbol{\alpha}'_2 \mathbf{x}, \dots, Z_p = \boldsymbol{\alpha}'_p \mathbf{x})$, we place $(K + 1)$
27 knots on each Z -coordinate and use the power basis, which is the collection of the
functions $1, Z_i, (Z_i - t_{ik})_+^3$, ($k = 1, 2, \dots, (K + 1)$; $i = 1, 2, \dots, p$). Conditions for
29 linearity are imposed similarly on each of the p coordinates.

4.2. The algorithm

31 Different steps of the IFS algorithm are as follows:

Step 1. Indicator variables Y_1, Y_2, \dots, Y_J are defined for J competing classes.

33 **Step 2.** Direction vectors $\boldsymbol{\alpha}_i$ ($\|\boldsymbol{\alpha}_i\| = 1$) ($i = 1, 2, \dots, p$) are initialized and the
linear features (combinations) Z_1, Z_2, \dots, Z_p are computed.

$$35 \quad Z_i = \boldsymbol{\alpha}'_i \mathbf{x}, \quad i = 1, 2, \dots, p. \quad (8)$$

1 **Step 3.** For each Z_i ($i = 1, 2, \dots, p$), K knots $t_{i_1} < t_{i_2} < \dots < t_{i_K}$ are placed in
 2 the range $[\min_n(Z_{i_n}), \max_n(Z_{i_n})]$. Basis functions are computed (as described in
 3 Sec. 4.1).

4 **Step 4.** Y_j is regressed on basis functions to get the initial estimates for
 5 $\phi_{j0}, \phi_{j1}, \dots, \phi_{jp}$, ($j = 1, 2, \dots, J$). Posterior probability estimates ($\hat{Y}_{j_n} = \psi_j(\mathbf{x}_n)$)
 6 and residuals ($r_{j_n} = Y_{j_n} - \psi_j(\mathbf{x}_n)$), ($j = 1, 2, \dots, J$, $n = 1, 2, \dots, N$) are computed
 7 to find the residual sum of squares (RSS) = $\sum_{j=1}^J \sum_{n=1}^N r_{j_n}^2$.

8 **Step 5.** Backfitting is used to readjust Z_1, Z_2, \dots, Z_p . At any stage α_1 is adjusted
 9 by a factor δ_1 , where δ_1 is the least squares solution obtained by minimizing

$$\text{RSS} \simeq \sum_{j=1}^J \sum_{n=1}^N [r_{j_n} - \delta_1 \eta_{j_n}^{(1)}]^2, \quad (9)$$

11 where $\eta_{j_n}^{(1)} = \frac{\partial \phi_{j1}}{\partial \alpha_1} |_{\mathbf{x}=\mathbf{x}_n}$ is computed at the current estimates of α_1 and ϕ_{j1} . This
 12 new estimate of α_1 is normalized, and Z_1 is recomputed. Knots are placed on
 13 the range of Z_1 to find the new basis functions. $Y_j - \sum_{i \neq 1} \phi_{ji}(Z_i)$ are regressed
 14 on those basis functions to estimate ϕ_{j0} and ϕ_{j1} ($j = 1, 2, \dots, J$). However, these
 15 adjustments are made only when this new estimate of α_1 reduces the current value
 16 of RSS, and in that case, probability estimates and residuals are also readjusted.
 17 Otherwise, all the estimates, knots and basis functions are kept unchanged. After
 18 α_1 , we proceed in the similar way to adjust α_2 and all other directions. Thus,
 19 Z_1, Z_2, \dots, Z_p are modified one by one. This step is repeated until no significant
 20 improvement is observed in RSS.

21 **Step 6.** Y_j is regressed on current basis functions to get new estimates for
 22 $\phi_{j0}, \phi_{j1}, \dots, \phi_{jp}$, ($j = 1, 2, \dots, J$). \hat{Y}_{j_n} and r_{j_n} ($j = 1, 2, \dots, J$; $n = 1, 2, \dots, N$)
 23 are adjusted accordingly.

24 The last two steps (Steps 5 and 6) are repeated until convergence is achieved
 25 (relative reduction in RSS is very small over a number of consecutive iterations).
 26 A detailed version of this algorithm is given in the Appendix.

27 Like CUS (see Ref. 3), IFS puts no restriction on ψ_j ($j = 1, 2, \dots, J$) to ensure
 28 that they are in $[0, 1]$. Imposing positivity restrictions by any manner results in much
 29 more complicated but not necessarily better classification (see Ref. 27). Due to spe-
 30 cial type of effect of bias-variance decomposition on misclassification rates (see e.g.
 31 Ref. 17), IFS leads to fairly good results in spite of having posterior estimates out-
 32 side the $[0, 1]$ range. However, inclusion of the intercept term in the set of basis func-
 33 tions, ϕ_{j0} ($j = 1, 2, \dots, J$) guarantees the additivity constraint ($\sum_{j=1}^J \hat{\psi}_j(\mathbf{x}_n) = 1$,
 $\forall n = 1, 2, \dots, N$).

35 4.3. The backward deletion procedure: Selection of optimum 36 dimensionality

37 As we have mentioned before, IFS algorithm starts with a reasonably large num-
 38 ber of initial knots. After fitting this full model, backward deletion is used to

1 attain models with smaller dimensionalities. At each step, we consider each of the
 2 undeleted knots as possible candidates for deletion, and delete that one which leads
 3 to least increase in the training set RSS. To maintain linearity at the left tail, linear
 4 basis functions cannot be considered as a candidate for deletion until all the knots
 5 corresponding to that variable have been deleted. During the deletion process, some
 6 restrictions are imposed to maintain the linearity of the fitted functions beyond
 7 the rightmost undeleted knots. If $t_{i_{(K+1)}}$, the rightmost knot on Z_i , is deleted,
 8 the condition of linearity to the right of t_{i_K} can be imposed by adding another
 9 constraint $\sum_{k=1}^K \beta_{i_k} = 0$. One needs to solve the least squares problem under
 10 one or two such constraints which are to be satisfied for each of the Z -variables.
 11 The whole deletion procedure is done by modified Gaussian sweep algorithm
 12 (see e.g. Ref. 7).

13 Backward deletion generates a sequence of nested models indexed by the dimen-
 14 sionality of the constrained space. The problem then reduces to selection of an
 15 optimum one under some criterion. As the optimum dimensionality depends on
 16 the problem itself, it is desirable to find it using the training set observations. If
 17 resubstitution error rate is used as a criterion, naturally one of the larger models
 18 would be selected which may not be the best one for classifying future observa-
 19 tions. Therefore, to arrive at a parsimonious model, we adopt the cost complexity
 20 criterion. Cost complexity for a model with dimension i is defined as

$$21 \quad R_\gamma(i) = R(i) + \gamma i, \quad i = 0, 1, 2, \dots \quad (10)$$

22 where $R(i)$ is the training set misclassification rate for the corresponding model,
 23 and γ is the cost complexity parameter. Clearly, $\gamma = 0$ leads to resubstitution
 24 error rate criterion when the cost complexity gets minimized by one of the larger
 25 models. This model remains optimum up to a certain positive value of γ after which
 26 another smaller model turns out to be the cost minimizer. In the process, a number
 27 of intervals and the corresponding minimizing models are obtained, and the number
 28 of competitive models usually gets reduced. Suppose we get m such intervals for
 29 $\gamma \{(\gamma_{t-1}, \gamma_t); t = 1, 2, \dots, m\}$, and the corresponding models are $M_1 \succ M_2 \succ \dots \succ$
 30 M_m . The problem of selection of final model then reduces to selection of an ideal
 31 cost parameter for which we use cross-validation.

32 In V -fold cross-validation (see e.g. Refs. 35 and 41), stratified random sampling
 33 is carried out to divide the whole training set into V groups L_1, L_2, \dots, L_V of
 34 sizes as nearly equal as possible. Observations belonging to different classes are
 35 used as different strata. Leaving one group L_r ($r = 1, 2, \dots, V$) at a time as a
 36 hold-out sample, IFS is used on the remaining observations. We start with the
 37 same number of initial knots per variable that was used with the whole training
 38 set, and then backward deletion is carried out in the same way. In the process,
 39 a nested sequence of models (as discussed above) is generated, and each time the
 40 resubstitution error rate is computed. The cost function is then minimized over
 41 these models for different γ , more specifically for $\gamma_t^* = \sqrt{\gamma_{t-1}\gamma_t}$, ($t = 1, 2, \dots, m$).
 Let $M_t^{(r)}$ be the models which minimize the cost functions $R_{\gamma_t^*}$ ($t = 1, 2, \dots, m$).

1 These models are then used to classify the observations in the hold out sample (L_r),
 2 and in each case the number of misclassifications is computed. We repeat the same
 3 procedure over the V groups to estimate the misclassification error rates (Δ) for
 4 different γ_t^* .

$$5 \quad \widehat{\Delta}(\gamma_t^*) = \frac{1}{N} \sum_{r=1}^V \sum_{\{n: (\mathbf{x}_n, c_n) \in L_r\}} I(d_t^{(r)}(\mathbf{x}_n) \neq c_n), \quad t = 1, 2, \dots, m, \quad (11)$$

6 where $d_t^{(r)}$ is decision rule obtained from the model $M_t^{(r)}$. The value of γ_t^* that
 7 minimizes $\widehat{\Delta}$, is used as the ideal value of the cost complexity parameter, and
 8 the corresponding model is selected as the final model. In our algorithm, we use
 9 $V = 10$ for cross-validation. Further details on cost-complexity pruning is available
 10 in Ref. 5.

11 **4.4. Number and placement of initial knots: Selection of initial 12 features**

13 Number of initial knots should be chosen carefully. To explain local patterns of the
 14 measurement space, this number should be reasonably large. From our simulation
 15 study, we feel that the result is not too sensitive as long as reasonably large number
 16 of knots are used. Our experience suggests that 10–12 knots per variable are enough
 17 for a moderately large sample size. Cross-validation can also be used to find this
 18 number. However, one must be careful when working with a large number of knots.
 19 Not only it leads to higher variability in the model selected but may also lead to
 20 a matrix $Z'Z$ which cannot be stably inverted. This problem of singularity can be
 21 tackled by carefully placing the knots. Equispaced knots in the range works well in
 22 many cases but performs poorly when Z_i s have large gaps between some consecutive
 23 observations. Knot placement based on order statistics seems to be a better choice.

24 Initial direction vectors α_i^0 ($i = 1, 2, \dots, p$) have to be specified as well. From
 25 our empirical experience we feel that instead of starting randomly, it is usually
 26 better to start with the linear features that maximize the linear separation among
 27 the classes. These direction vectors are the eigen vectors of $W^{-1}(W + B)$, where B
 28 and W stand for between class and within class sum of square matrices. Of course
 29 one can also start with the *CUS* model where α_i^0 s are taken as the unit vectors
 30 along the co-ordinate axes (i.e. the original variables as the initial linear features).
 31 In the results reported in Sec. 5, IFS (features) and IFS (variables) represent the
 32 results obtained using these two different starting points, respectively.

33 **4.5. Convergence issues**

34 There are essentially three minimization procedures in IFS. Minimization over the
 35 linear features (the α parameters) which is done iteratively, minimization over the
 36 class of additive spline functions which is done by estimating the β parameters
 37 using least squares and minimization over the number of linear features which is

1 done by cost-complexity cross-validation. While convergence of the entire algorithm
 2 requires rather involved mathematical exercise, results regarding each of the three
 3 minimization are available in the literature which provides justification for the
 4 effectiveness of the algorithm.

5 The minimization over the α parameters is done by backfitting which was effec-
 6 tively used by Breiman and Friedman.⁶ For estimating the additive transforma-
 7 tions of several explanatory variables, they considered stepwise minimization of the
 8 error sum of squares involving the transformation function of only one variable at
 9 a time, and repeating the procedure until convergence. The advantage is that at
 10 each step, the error sum of squares is reduced by least squares which guarantees
 11 convergence. However there may not be a unique minimum for this part of the
 12 optimization.

13 Given the linear features, the IFS procedure is essentially the CUS procedure
 14 on the transformed feature space instead of the original measurement space. Con-
 15 vergence of IFS has been proved in Ref. 2. The proof relies on the universal approx-
 16 imation property of cubic splines.

17 Finally the optimality of the cost-complexity cross-validation has been discussed
 18 in detail by Breiman *et al.*⁵ We have also studied the convergence empirically with
 19 a simulated dataset in Sec. 5.1.

4.6. Problem of multiple local minima: Use of simulated annealing

21 IFS may suffer from similar problems like neural networks because of the presence
 22 of possibly numerous local minima in the transformed feature space. Therefore,
 23 it may be worthwhile to start with different starting points or to use simulated
 24 annealing. We opted for the second one to keep the algorithm automatic. During
 25 feature extraction, each time we compute two competitive direction vectors, one
 26 as discussed in Sec. 4.2 and another by simulation from a uniform distribution
 27 (over a unit hyper-sphere with the current value of α_i as origin). After normal-
 28 ization, we chose the one which has the smaller training set RSS as the new esti-
 29 mate α_i^* . If this new estimate reduces the current value of RSS, Z_i is modified.
 30 Otherwise, we draw a random number u from $U(0, 1)$ and adjust Z_i only when
 31 $u \leq \exp\left(-\frac{\text{RSS}^* - \text{RSS}}{\text{RSS} \times C_m}\right)$, where RSS^* is the residual sum of squares corresponding to
 32 α_i^* . This step helps to overcome the problem of getting stuck at bad local minima.
 33 C_m is a constant that depends on the number of iteration (m). Generally, a large
 34 number is chosen as C_0 to make almost all the directions acceptable in the begin-
 35 ning. It is then decreased gradually over iterations to attain convergence. Several
 36 options are available for determining these constants (see Ref. 28), but we use a
 37 simple one. We start with $C_0 = \frac{100}{p}$, where p is the number of measurement vari-
 38 ables, and then use $C_m = U(0.8, 0.99) \times C_{m-1}$ ($m = 1, 2, \dots$) to modify it. During
 39 modification of direction vectors α and features ϕ , at any stage, we always store
 40 the linear combinations Z_1, Z_2, \dots, Z_p for which training set RSS is minimum up to
 41 that stage.

1 5. Experimental Results

3 In this section, we use some simulated and benchmark data sets to illustrate the
4 usefulness of the proposed method. In Tables 1 and 2, we report the misclassifi-
5 cation rates of IFS algorithm starting with two different sets of direction vectors
6 as mentioned in Sec. 4.4 and denote them by IFS (variables) and IFS (features),
7 respectively. For assessing the accuracy of IFS, we compare their performances with
8 a number of existing parametric and nonparametric classifiers. Results of linear dis-
9 criminant analysis (LDA), quadratic discriminant analysis (QDA), CUS, method of
10 successive projections, neural networks (with sigmoidal transformation function),
11 radial basis function networks (RBF), kernel discriminant analysis, nearest neighbor
12 classifiers and support vector machines (SVM) are reported to facilitate the com-
13 parison. Error rates for classification tree method are also reported in benchmark
14 data sets. For simulated data sets, we compute the optimal error rates and present
15 them in Table 1. These optimal error rates represent the misclassification rate of
16 the optimal Bayes rule applied to specific data sets.

17 For kernel discriminant analysis (see e.g. Ref. 22), we standardized the obser-
18 vations of a class by using the usual moment based estimate of class dispersion
19 matrix and then used same bandwidth in all directions. Optimal bandwidth for a
20 class density estimate was selected by least square cross-validation (see e.g. Refs. 38
21 and 40). To find density estimates at a new data point, at first we standardized it
22 using the same standardization matrices used earlier to get the density estimate at
23 the standardized data points. Density estimates at the original data point can be
24 computed from that using a simple formula when the measurement vector under-
25 goes a linear transformation. For nearest neighbor classification (see e.g. Refs. 10
26 and 18), we used Mahalanobis distance (see e.g. Ref. 30) as the distance function,
27 which is equivalent to use Euclidean distance after standardization by estimated
28 pooled dispersion matrix.

29 For the neural networks (single hidden layer perceptrons with sigmoidal transfor-
30 mation) and the radial basis function (RBF) networks, we standardized the mea-
31 surement vectors before using training algorithms. For the perceptrons, we tried
32 both backpropagation (see e.g. Ref. 35) and Levenberg–Marquardt algorithms (see
33 e.g. Ref. 39) and opted for the latter because of its better performance both in
34 terms of computing time and misclassification rate. On each data set, we trained
35 the model for several choices of hidden nodes, and reported the best error rate that
36 we got in the test set. Similarly, we used various choices of bandwidth in the case of
37 radial basis function networks, and several choices of bandwidth and cost parameter
38 in the case of support vector machines, and reported their best performances on
39 test sets. Matlab programs available in Matlab neural network toolbox and Matlab
40 SVM toolbox were used for running these three algorithms. In the case of SVM,
41 for multiclass problems, voting was used to take the final decision after all pairwise
42 comparisons. We used the S-Plus package for building the classification trees. Codes
43 for other classification algorithms were written in C-language.

1 **5.1. Results from the analysis on simulated data sets**

2 Here, we consider three simulated examples on two-class classification problems.
 3 The first two of these data sets contain two-dimensional measurement vectors, while
 4 in Example 3, we consider a higher dimensional problem. Since training sets for
 5 real problems are usually not very large (because of the cost involving generation
 6 of samples), we use relatively smaller training sets in our simulations. However, to
 7 obtain reliable estimates for misclassification rates, much larger test sets are used.
 8 For each of these examples, we generate 500 observations for training and 3000
 9 for the test sets taking equal number of observations from the two classes. Each
 10 experiment is carried out ten times, and average misclassification rates of different
 11 classifiers over these ten simulation runs are reported in Table 1 along with their
 12 corresponding standard errors.

13 **Example 1.** In this example, each class is an equal mixture of two bivariate normal
 14 distributions differing only in their location parameters. The population density
 15 functions are given by

$$\begin{aligned} p(\mathbf{x}|1) &= 1/2[N_2(1, 1, 0.6^2, 0.6^2, 0) + N_2(-1, -1, 0.6^2, 0.6^2, 0)] \\ p(\mathbf{x}|2) &= 1/2[N_2(1, -1, 0.6^2, 0.6^2, 0) + N_2(-1, 1, 0.6^2, 0.6^2, 0)] \end{aligned} \quad (12)$$

17 Clearly, $X_1X_2 = 0$ is the optimal class boundary for this problem. Therefore, the
 18 ideal linear combinations (linear features) are of the form $Z_1 = \alpha_1X_1 + \alpha_2X_2$ and
 19 $Z_2 = \alpha_1X_1 - \alpha_2X_2$ ($\alpha_1^2 + \alpha_2^2 = 1$), whereas the corresponding features (trans-
 20 formation functions) are quadratic in nature. It is evident from Table 1 that IFS
 21 performed quite well in this example, and in fact, it could estimate the ideal linear
 22 combinations and transformation functions very accurately in all ten simulation
 23 runs, one of which is chosen below for illustration.

24 In this example, linear combinations Z_1 and Z_2 were estimated as $0.705X_1 +$
 25 $0.709X_2$ and $0.676X_2 - 0.737X_1$. Due to additivity constraint, as expected, we got
 26 $\hat{\phi}_{1i} = -\hat{\phi}_{2i}$ (for $i = 1, 2$) and $\hat{\phi}_{10} + \hat{\phi}_{20} = 1$. Moreover, the desired quadratic nature
 27 of the plotted functions $\hat{\phi}_{11}$ and $\hat{\phi}_{12}$ (see Fig. 1) clearly suggests the success of IFS
 28 in estimating the ideal functions as well. Though instead of actual values of X_1X_2 ,
 29 we used only the indicator variables, IFS was still able to detect the proper linear
 30 combinations and the corresponding nonlinear features (transformation functions).
 31 Hence the estimated boundaries in Fig. 2(a) also turned out to be quite satisfactory.

32 Table 1 shows that in this example, performance of LDA was not satisfactory at
 33 all but error rates for QDA came closest to the optimal error rates. Classification
 34 methods like CUS and successive projections, which are based on additive regression
 35 did not work at all. However, IFS and other nonparametric classifiers could nearly
 36 match the performance of QDA. As the optimal class boundary ($X_1X_2 = 0$) is not
 37 an additive function of the measurement variables, in a few cases IFS got stuck to
 38 bad local minima when original variables are used as the initial features. Simulated
 39 annealing helped in such situations. In this example, successive projections could
 not improve upon CUS, but IFS did a remarkable job for this purpose.

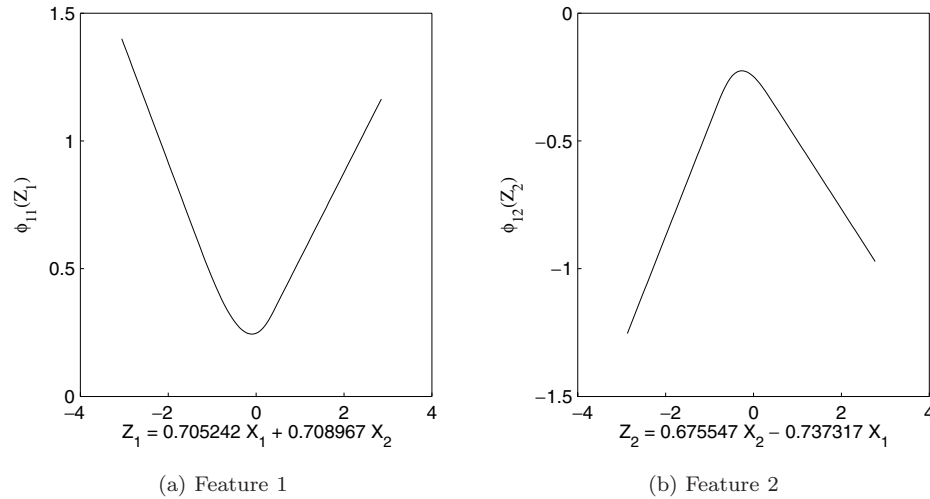


Fig. 1. Estimated features in Example 1.

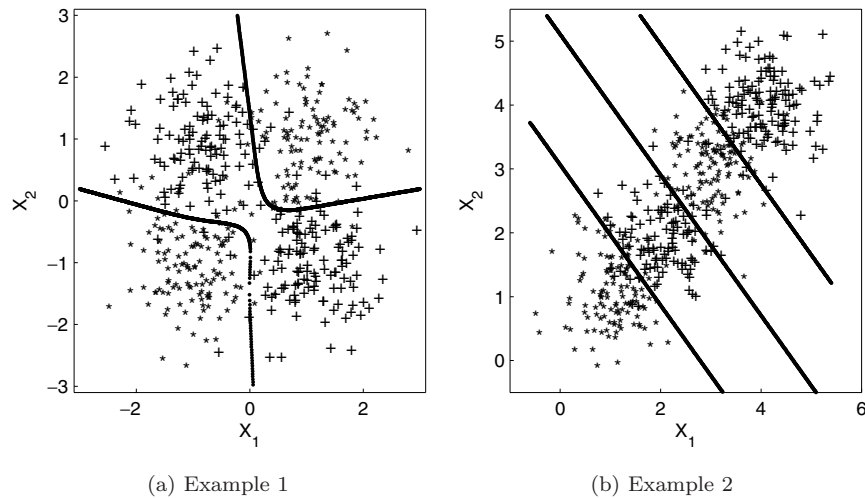


Fig. 2. Estimated class boundaries in Examples 1 and 2.

- 1 **Example 2.** Like the previous example, here also competing classes are equal mix-
 2 tures of bivariate normal populations but overlaps between the competing classes
 3 occur in a different manner. Population density functions of these two populations
 4 are given by

$$\begin{aligned}
 p(\mathbf{x}|1) &= 1/2[N_2(1, 1, 0.5^2, 0.5^2, 0) + N_2(3, 3, 0.5^2, 0.5^2, 0)] \\
 p(\mathbf{x}|2) &= 1/2[N_2(2, 2, 0.5^2, 0.5^2, 0) + N_2(4, 4, 0.5^2, 0.5^2, 0)]
 \end{aligned}
 \tag{13}$$

- 5
 6 Since the true class boundary is highly nonlinear in this example, neither LDA nor
 7 QDA could manage to estimate it properly (see Table 1). CUS and the method of

Table 1. Average misclassification rates and their standard errors (in percentage) in simulated examples.

	Example-1	Example-2	Example-3
Optimal	9.04 (0.15)	12.08 (0.16)	10.74 (0.14)
LDA	49.18 (0.50)	45.23 (0.57)	49.98 (0.33)
QDA	9.24 (0.11)	45.44 (0.54)	11.61 (0.17)
CUS	48.23 (0.83)	18.35 (0.41)	14.30 (0.27)
Succ. Proj.	44.55 (1.83)	18.54 (0.40)	14.02 (0.28)
Nearest neighbor	9.62 (0.13)	13.22 (0.19)	16.49 (0.20)
Kernel	9.80 (0.14)	12.86 (0.23)	12.43 (0.16)
SV-Machine	9.67 (0.14)	12.41 (0.18)	13.04 (0.18)
Neural network	9.74 (0.15)	12.58 (0.21)	12.48 (0.19)
RBF network	9.46 (0.16)	12.33 (0.17)	12.92 (0.17)
IFS (variables)	9.76 (0.17)	12.58 (0.18)	12.31 (0.18)
IFS (features)	9.66 (0.20)	12.38 (0.18)	11.88 (0.21)

1 successive projections showed much better performance, yet the misclassification
 2 rates were far from optimum. However, other nonparametric classifiers achieved
 3 reasonably low misclassification rates. Among these classifiers, nearest neighbor
 4 method had slightly higher error rate but performances of other classifiers were
 5 fairly similar. It is quite evident from Fig. 2(b) that IFS could identify the optimal
 6 class boundary. Though successive projections failed to improve the performance
 7 of CUS, IFS was very successful in this problem again.

8 **Example 3.** Next we consider a higher dimensional classification problem between
 9 two populations, where the first two co-ordinate variables in the two classes are
 10 distributed as

$$11 \quad p(\mathbf{x}|1) = N_2(0, 0, 1, 1, 1/2) \quad p(\mathbf{x}|2) = U[-5, 5] \times U[-5, 5] \quad (14)$$

12 Other three variables are independent and identically distributed as $N(0, 1)$ in both
 13 the populations. These variables do not contain any separability information, and
 14 they are used simply to add noise.

15 Clearly, in this example, the true class boundary is of the form $X_1^2 + X_2^2 -$
 16 $X_1X_2 = C_0$ for some constant C_0 . This can be rewritten (not uniquely) as $A\{(X_1 -$
 17 $0.5X_2)\}^2 + BX_2^2 = C$, where A, B and C are appropriate constants. Thus, we see
 18 that only two linear features contain the information about class separability. When
 19 we ran IFS algorithm on different sets of 500 observations generated from these two-
 20 classes, the final model contained exactly two linear features in most of the cases.
 21 Figure 3 shows that IFS could obtain very appropriate estimates of linear features
 22 and that of the optimal class boundary in this noisy example. Due to quadratic
 23 nature of the true class boundary, as expected, QDA led to the best performance in
 24 this example. IFS performed better than all other parametric and nonparametric
 25 classifiers and the corresponding error rates were close to that of QDA.

26 Since the optimal error rate is known for a simulated problem, we chose to per-
 27 form an empirical study. In Example 1, we trained IFS on training sets of varying

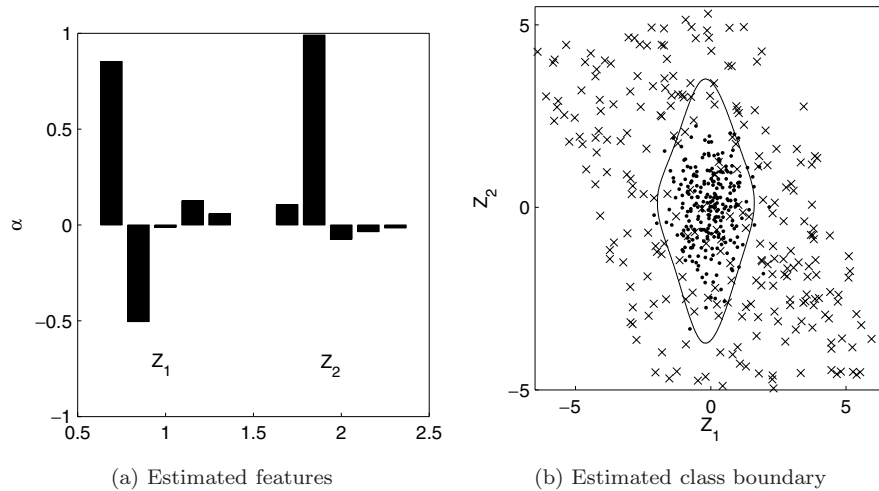


Fig. 3. Estimated features and class boundaries in Example 3.

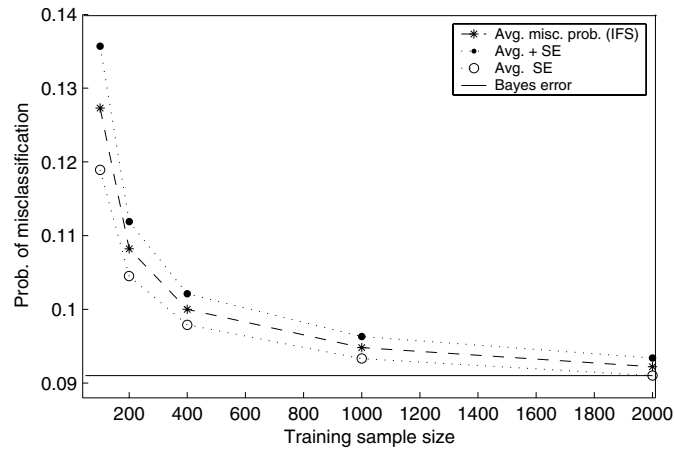


Fig. 4. Misclassification rates of IFS for different sample sizes.

1 sizes consisting of 100, 200, 400, 1000 and 2000 cases, and computed the test set
 2 errors based on fixed size of 3000 test cases. Figure 4 shows the average misclassifi-
 3 cation rate over these ten simulation runs and the corresponding standard error for
 4 each training sample size. This figure clearly indicates that the test set error rate
 5 of IFS converges quite quickly to some value very close to the optimal error rate.

5.2. Results from the analysis on benchmark data sets

7 We analyze five benchmark data sets for further illustration of the proposed method.
 8 All these data sets has previously been used by many other authors (see e.g. Refs. 3,
 9 9, 23, 24, 34 and 36), who evaluated the performance of many classifiers on these

16 *A. K. Ghosh & S. Bose*

1 data sets. Along with the error rates of all parametric and nonparametric classi-
 3 fiers that we considered in Sec. 5.1, here we report the misclassification rates for
 5 classification tree methods as well. Throughout this section, sample proportions of
 7 different classes are used as their prior probabilities. The data sets that we use in
 this section have specific training and test samples. The sizes of the training and
 test sets for each data set have been reported in Table 2, which also shows the test
 set error rates for different classification methods. A brief description of these data
 sets is given below.

9 **Synthetic data.** It consists of bivariate observations from two competing classes,
 11 each of which is an equal mixture of two bivariate normal populations. All these
 13 normal populations have the same dispersion matrix but different location param-
 eters. These parameters were chosen to have a Bayes risk of 8.0%. A scatter plot
 of this data set is given in Fig. 4. This data set is available at CMU data archive
 (<http://lib.stat.cmu.edu>).

15 We have used two data sets related to vowel recognition problem and refer them
 as vowel data-1 and vowel data-2.

17 **Vowel data-1.** This data was created by Peterson and Barney³³ by a spectro-
 19 graphic analysis of vowels in words formed by “h” followed by a vowel and then
 followed by “d”. There were 67 persons who spoke different words and the two
 lowest resonant frequencies of a speaker’s vocal track were noted for ten different
 21 vowels. A scatter plot of this data set is given in Fig. 6, where the numbers represent
 the labels of different classes (“0” represents the tenth class).

23 **Image data.** It originally contains 19 measurements on each image of one of seven
 different objects : brickface, sky, foliage, cement, window, path and grass. This

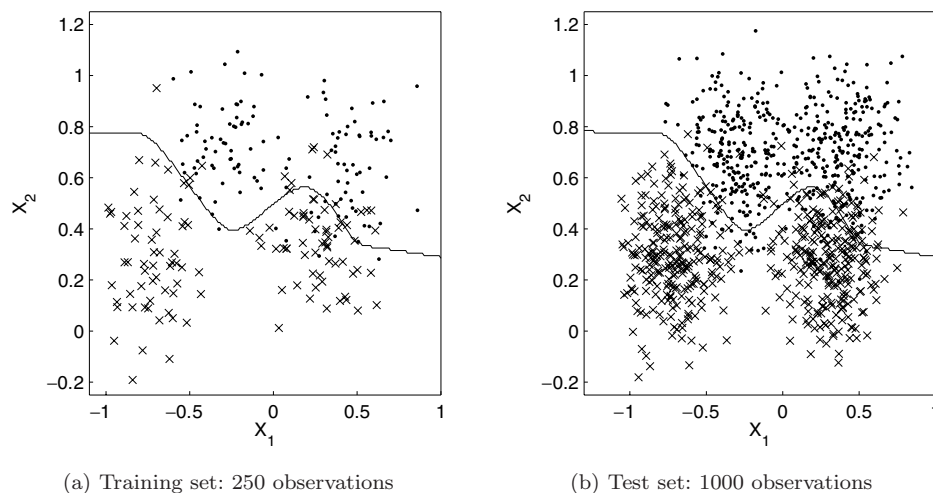


Fig. 5. Estimated class boundaries for synthetic data.

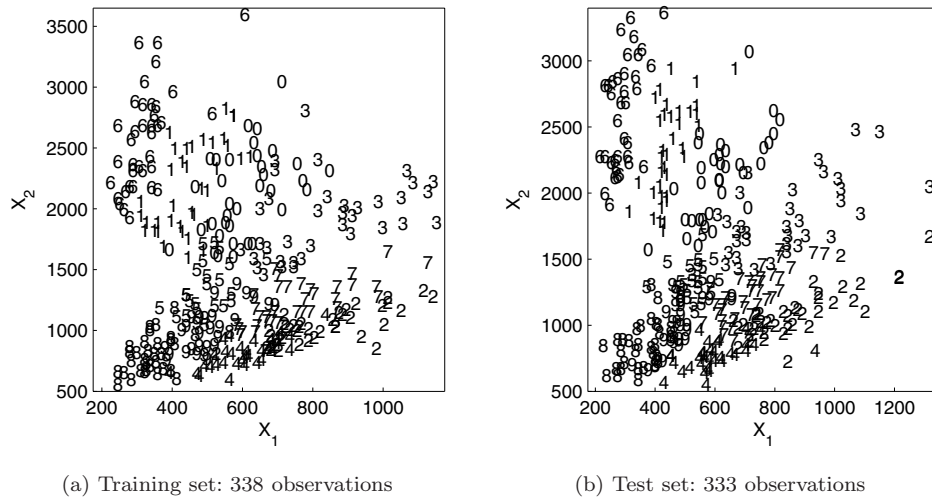


Fig. 6. Scatter plots for vowel data-1.

1 data set and the description of the variables are available at UCI machine learning
 2 repository (<http://www.ics.uci.edu/~mllearn>). The value of the variable ‘region
 3 pixel count’ is ‘9’ for all observations. For other two variables, “short line density-5”
 4 and “short line density-2”, almost 95% of the values are zero. We did not consider
 5 these three variables in our study. There are some variables in the data set which are
 6 linear or nonlinear functions of three other variables: R (“raw red mean”), B (“raw
 7 blue mean”) and G (“raw green mean”). We have deleted those variables too and
 8 carried out our analysis using the remaining nine variables (region-centroid column,
 9 region-centroid-row, vedge-mean, vedge-sd, hedge-mean, hedge-sd, rawread-mean,
 10 rawblue-mean and rawgreen-mean).

11 **Satimage data.** The sample database was generated by taking a small section
 12 from the original Landsat Multi-Spectral Scanner Image Data purchased from
 13 NASA by the Australian Centre for Remote Sensing. The database consists of
 14 the multispectral values of pixels (coded as numbers) in 3×3 neighborhoods in
 15 a satellite image, and the classification associated with the central pixel in each
 16 neighborhood. However, for our analysis we used only four spectral values related
 17 to central pixel to classify it into one of six different classes : red soil, cotton crop,
 18 grey soil, damp grey soil, soil with vegetation stubble and very damp grey soil. The
 19 data set and its description is available at UCI machine learning repository.

20 **Vowel data-2.** This data set is also available at UCI machine learning repository.
 21 Like vowel data-1, it is also related to a vowel recognition problem, where ten
 22 measurements are taken for each observation coming from any of 11 classes. The
 23 description of the measurement variables can be found in Ref. 23. There are 528
 24 observations in the training set and 462 observations in the test set, which are
 25 equally distributed among 11 competing classes.

Table 2. Misclassification rates (in percentage) for benchmark data sets.

	Synthetic	Vowel-1	Satimage	Image	Vowel-2
Dimension	2	2	4	9	10
No. of classes	2	10	6	7	11
Train. sample size	250	338	4435	210	528
Test sample size	1000	333	2000	2100	462
LDA	10.80	25.23	19.30	11.14	55.63
QDA	10.20	19.82	15.65	14.62	52.81
CUS	8.50	24.32	19.55	10.10	51.30
Succ. Proj.	8.50	18.92	16.90	7.48	45.67
Nearest Neighbor	11.70	17.72	15.35	8.24	46.75
Classification Tree	10.10	23.72	19.50	12.57	56.28
Kernel	9.30	18.32	15.30	15.71	62.12
SV-Machine	9.40	21.02	14.80	11.95	46.54
Neural network	9.40	18.62	16.95	9.95	48.92
RBF network	9.40	20.42	15.70	14.00	49.78
IFS (variables)	9.10	20.72	15.80	6.76	44.37
IFS (features)	9.60	23.42	15.40	6.28	43.29

1 Table 2 reveals that IFS performed quite well in all these data sets. On syn-
 2 thetic data, CUS and successive projection methods had the best error rates but IFS
 3 and all other nonparametric methods except classification tree and nearest neighbor
 4 classifier worked well. Class boundary estimated by IFS (features) is shown in Fig. 5.
 5 On vowel data-1, nearest neighbor method led to the best error rate. Kernel discrim-
 6 inant analysis and successive projection methods also had reasonably lower error
 7 rate. Performances of IFS and other nonparametric methods were fairly satisfactory.
 8 On satellite image data (satimage data), SVM yielded the lowest misclassification
 9 rate but the error rates for QDA, kernel discriminant analysis, nearest neighbor,
 10 RBF neural network and IFS were very close to that. Performance of IFS was much
 11 better than CUS and successive projection methods. On image segmentation data
 12 (image data), IFS performed extremely well, and the test set misclassification rates
 13 for the two different starting points were much lower compared to other parametric
 14 and nonparametric classification techniques. Given that the test sample size was
 15 quite large, the difference between the error rate of IFS and that of the other clas-
 16 sifiers was found to be statistically significant. On vowel data-2, once again IFS led
 17 to best error rates among the parametric and nonparametric classifiers considered
 18 here. Successive projections could reduce the misclassification rate of CUS in this
 19 example quite a bit, but IFS did even better. This data set was used extensively
 20 by Hastie and Tibshirani and Hastie *et al.*,²³ where they reported results of many
 21 other classification methods. IFS performed better than every other method except
 22 for FDA-MARS (deg-2) which had a slightly lower error rate. Overall the error rate
 23 of IFS was either the lowest among the methods compared or was very close to the
 24 lowest error rate for a given problem. IFS also had the lowest average test set error
 25 across all the examples.

26 Since different classification methods were run in different platforms in our
 27 experiments, it was not possible to compare their CPU times. Perhaps with the

Table 3. Average number of iterations and CPU times taken by IFS algorithms for training the network on different data sets (figures in top and bottom rows are for IFS(variables) and IFS(features), respectively). These averages are taken over 11 runs including ten times for cross-validation.

	Example-1	Example-2	Example-3	Synthetic	Vowel-1	Satimage	Image	Vowel-2
No. of iterations	15.9	12.8	22.4	11.7	22.9	24.7	94.2	83.5
	13.5	12.7	33.7	9.3	10.5	30.9	91.5	67.6
CPU time (secs.)	0.74	0.63	2.92	0.28	2.05	86.05	28.53	133.80
	0.65	0.62	4.33	0.23	1.07	108.87	28.38	104.39

1 evolution of faster computers, the training time is not as important as it used to
 2 be. However, IFS (implemented in C-programming language) works reasonably fast.
 3 Since it is an iterative algorithm, it is difficult to find its computational complexity.
 4 If we consider the dimension d and the number of class J to be finite (which is
 5 usually the case) for each iteration IFS requires $O(K^2n)$ calculations, where K is
 6 the number of knots placed on the range of each variable and n is the training
 7 sample size. The number of knots should increase with the sample size but at an
 8 extremely slow rate. Hence, the computational cost per iteration is only marginally
 9 higher than $O(n)$. Table 3 below shows the CPU time (on a Pentium-4 machine)
 10 and average number of iterations required by IFS algorithm on different data sets.
 11 For instance, on vowel data-2 it took around 20 minutes to train the model 11 times
 12 including ten times for cross-validation. However, on the same platform, for a single
 13 run of neural network (a two-layer perceptron with 20 hidden nodes), MATLAB
 14 took almost half an hour (for 100 iterations) for training. For complex problems,
 15 generally a large number of hidden nodes is required and hence it takes a long time
 16 to train the network. But IFS deals with fewer number of linear features which
 17 helps to reduce the computing cost substantially.

6. Conclusion

19 In real life classification problems, the class boundaries are more complex rather
 20 than being linear or quadratic. In such cases, traditional methods of discriminant
 21 analysis often turn out to be inadequate. This article presents a nonparametric
 22 method (IFS) for discriminant analysis, which is capable to approximate these
 23 class boundaries using a model similar to the projection pursuit regression model. It
 24 automatically tries to identify some linear combination of the original measurement
 25 variables and fits an additive model based on regression splines using those combi-
 26 nations. Model selection (selection of the ideal number of linear combinations and
 27 the number of knots for fitting splines) is done using cross-validation.

28 IFS can be viewed as a generalized version of the single hidden layer perceptron
 29 model, which uses sigmoidal or other known transformation functions at the hidden
 30 layer. On the other hand, instead of using any fixed transformation function, IFS
 31 uses cubic splines to estimate the transformation functions from the data itself
 which makes it more flexible.

20 *A. K. Ghosh & S. Bose*

1 The use of backfitting for convergence makes the algorithm fast. This enables
 2 use of cross-validation for model selection which requires a number of training on
 3 different partitions of the data. Model selection by backward deletion allows for
 4 very simple models.

5 IFS can also be viewed as a modification of CUS (see Ref. 3), which uses an addi-
 6 tive model based on the original measurement variables. The experiments reported
 7 in Sec. 5 clearly show the effectiveness of IFS over CUS in terms of lower mis-
 8 classification rates in various classification problems. Another advantage of being
 9 a generalization of the CUS procedure is that IFS does not have to use random
 10 starting points. If one starts with α_i^0 which are unit vectors along co-ordinate axes
 11 (i.e. original variables as initial linear features), the initial result will be the result
 12 of CUS which usually yields reasonable misclassification rates. The subsequent iter-
 13 ations can only improve the results further. Thus one can expect convergence to a
 14 better solution.

15 The experiments also show that IFS has a definite edge over perceptron models
 16 both in terms of misclassification rate and also in visualizing the class separability.
 17 It could find out the ideal direction vectors and nonlinear features in all examples
 18 that we have tried including the ones reported here. IFS also fared well compared
 19 to many other nonparametric classification techniques. Finally the entire procedure
 20 is fully automatic. Only the initial number of knots has to be specified.

21 IFS can be adopted for regression problems also. Instead of indicator variables,
 22 the response variables can be used to estimate the regression surface by projection
 23 pursuit regression models. One of the major issues in projection pursuit regression
 24 is to find the proper linear combinations. IFS uses backfitting to estimate them
 25 iteratively and automatically. A similar work using smoothing splines is available
 26 in Ref. 36, which perhaps could be improved upon using a model selection procedure
 27 such as the one used in IFS.

Acknowledgment

29 We are thankful to an associate editor and four anonymous referees who carefully
 30 read an earlier version of the paper and provided us with several helpful comments.

Appendix

31 IFS Algorithm:-

||Initialization||

Define indicator variables Y_1, Y_2, \dots, Y_J for J competing classes.

Initialize the direction vectors $\alpha_1^0, \alpha_2^0, \dots, \alpha_p^0$ ($\|\alpha_i^0\| = 1 \forall i$) and compute the

linear combinations $Z_i^0 \leftarrow \alpha_i^{0'} \mathbf{x}$, $i = 1, 2, \dots, p$.

For each Z_i^0 , place the knots $t_{i1}^0, t_{i2}^0, \dots, t_{iK}^0$ on its range to define the basis functions

$$B_{i0}^0 \leftarrow 1, B_{ik}^0 \leftarrow Z_{i0}^0, B_{ik}^0 \leftarrow (Z_i^0 - t_{ik}^0)_+^3, \quad i = 1, 2, \dots, p; k = 1, 2, \dots, K.$$

Regress Y_j ($j = 1, 2, \dots, J$) on the basis functions to estimate $\phi_{j0}^0, \phi_{j1}^0, \dots, \phi_{jK}^0$.

Compute probability estimates $\hat{Y}_{j_n}^0 \leftarrow \phi_{j_0}^0 + \sum_{i=1}^p \phi_{j_i}^0(Z_{i_n}^0)$ and residuals
 $r_{j_n}^0 \leftarrow Y_{j_n} - \hat{Y}_{j_n}^0$, ($n = 1, 2, \dots, N$).
Count $\leftarrow 0$; Stop $\leftarrow 0$; $m \leftarrow 0$, $\text{RSS}_0 \leftarrow \sum_{j=1}^J \sum_{n=1}^N (r_{j_n}^0)^2$.
while(Stop=0) do
 ||Backfitting||
 $\text{RSS}_a \leftarrow \text{RSS}_0$
 for ($i = 1$ to p) do
 Use Taylor series approximation (up to linear term) of RSS about α_i^m

$$\text{RSS} \simeq \sum_{j=1}^J \sum_{n=1}^N [r_{j_n}^m - \delta_i' \xi_{j_n}]^2,$$
 where $\delta_i = \alpha_i - \alpha_i^m$ and $\xi_{j_n} = \frac{\partial \phi_{j_i}(\alpha_i' \mathbf{x})}{\partial \alpha_i} |_{\alpha_i = \alpha_i^m, \mathbf{x} = \mathbf{x}_n}$.
 Use least squares method to estimate δ_i .
 Compute $\alpha_i^{m*} \leftarrow \frac{\alpha_i^m + \delta_i}{\|\alpha_i^m + \delta_i\|}$ and $Z_i^{m*} \leftarrow \sigma(\alpha_i^{m*'} \mathbf{x})$.
 Place the knots $t_{i_k}^{m*}$ on Z_i^{m*} to compute basis functions $B_{i_k}^{m*}$, ($k = 0, 1, \dots, K$).
 Regress $Y_j - \sum_{k < i} \phi_{j_k}^{m+1}(Z_k^{m+1}) - \sum_{k > i} \phi_{j_k}^m(Z_k^m)$ on these basis functions to
 estimate $\phi_{j_0}^{m*}$ and $\phi_{j_i}^{m*}$, ($j = 1, 2, \dots, J$).
 Compute $\hat{Y}_{j_n}^{m*}$, $r_{j_n}^{m*}$, ($j = 1, 2, \dots, J$; $n = 1, 2, \dots, N$) and RSS^* accordingly.
 if($\text{RSS}^* \leq \text{RSS}_a$)
 $\text{RSS}_a \leftarrow \text{RSS}^*$, $\alpha_i^{m+1} \leftarrow \alpha_i^{m*}$; $Z_i^{m+1} \leftarrow Z_i^{m*}$.
 $B_{i_k}^{m+1} \leftarrow B_{i_k}^{m*}$, ($k = 1, 2, \dots, K$).
 $\phi_{j_i}^{m+1} \leftarrow \phi_{j_i}^{m*}$; $\phi_{j_0}^m \leftarrow \phi_{j_0}^{m*}$, ($j = 1, 2, \dots, J$).
 $\hat{Y}_{j_n}^m \leftarrow \hat{Y}_{j_n}^{m*}$; $r_{j_n}^m \leftarrow r_{j_n}^{m*}$, ($j = 1, 2, \dots, J$; $n = 1, 2, \dots, N$).
 end(**if**)
 if($\text{RSS}^* > \text{RSS}_a$)
 $\alpha_i^{m+1} \leftarrow \alpha_i^m$; $Z_i^{m+1} \leftarrow Z_i^m$; $\phi_{j_i}^{m+1} \leftarrow \phi_{j_i}^m$, ($j = 1, 2, \dots, J$).
 end(**if**)
 end(**for**)
 $\phi_{j_0}^{m+1} \leftarrow \phi_{j_0}^m$; $\hat{Y}_{j_n}^{m+1} \leftarrow \hat{Y}_{j_n}^m$; $r_{j_n}^{m+1} \leftarrow r_{j_n}^m$, ($j = 1, 2, \dots, J$; $n = 1, 2, \dots, N$).
 $\text{RSS}_1 \leftarrow \text{RSS}_a$; Reduction $\leftarrow \frac{\text{RSS}_1 - \text{RSS}_0}{\text{RSS}_0}$.
 if(Reduction ≤ 0.001)
 Regress Y_j ($j = 1, 2, \dots, J$) on the current basis functions to re-compute
 $\phi_{j_0}^{m+1}, \phi_{j_1}^{m+1}, \dots, \phi_{j_p}^{m+1}$.
 Re-adjust $\hat{Y}_{j_n}^{m+1}, r_{j_n}^{m+1}$, ($n = 1, 2, \dots, N$) and RSS_1 .
 Reduction $\leftarrow \frac{\text{RSS}_1 - \text{RSS}_0}{\text{RSS}_0}$.
 end(**if**)
 $m \leftarrow m + 1$; $\text{RSS}_0 \leftarrow \text{RSS}_1$.
 ||Termination||
 if(Reduction > 0.001) Count $\leftarrow 0$
 if(Reduction ≤ 0.001) Count $\leftarrow \text{Count} + 1$
 if((Reduction=0) or (Count=Count_stop)) Stop $\leftarrow 1$
end(**while**)

22 A. K. Ghosh & S. Bose

- 1 • Count_Stop is a user defined parameter for stopping criterion. In this article, we
 2 use Count_Stop = 5 and terminate the program if relative reduction in RSS is
 3 insignificant (i.e. Reduction < 0.001) over five consecutive iterations. However,
 4 the final result is not very sensitive to this parameter if any larger value is used.
 5 • For the sake of simplicity, the simulated annealing part has not been included
 6 in the above description, but one can easily incorporate it following the idea
 7 described in Sec. 4.6.

References

- 9 1. T. W. Anderson, *An Introduction to Multivariate Statistical Analysis* (Wiley, New
 10 York, 1984).
 11 2. S. Bose, A method for estimating nonlinear class boundaries in the classification
 12 problem and comparison with other existing methods, PhD Dissertation. (Dept. of
 13 Stat., Univ. of California, Berkeley, 1992).
 14 3. S. Bose, Classification using splines, *Comput. Stat. Data Anal.* **22** (1996) 505–525.
 15 4. S. Bose, Multilayer statistical classifier, *Comput. Stat. Data Anal.* **42** (2003)
 16 685–701.
 17 5. L. Breiman, J. H. Friedman, R. A. Olshen and C. J. Stone, *Classification and Regres-
 18 sion Trees* (Wadsworth & Brooks, Monterey, California, 1984).
 19 6. L. Breiman and J. H. Friedman, Estimating optimal transformations for multiple
 20 regression and correlation (with discussion), *J. Amer. Stat. Assoc.* **80** (1985) 580–
 21 619.
 22 7. L. Breiman, Fitting additive models to regression data: diagnostics and alternating
 23 views. *Comput. Stat. Data Anal.* **15** (1993) 13–46.
 24 8. B. Cheng and D. M. Titterton, Neural networks: a review from a statistical per-
 25 spective (with discussion), *Stat. Sci.* **9** (1994) 2–54.
 26 9. C. A. Cooley and S. N. MacEachern, Classification via kernel product estimators,
 27 *Biometrika* **85** (1998) 823–833.
 28 10. T. M. Cover and P. E. Hart, Nearest neighbour pattern classification, *IEEE Trans.
 29 Inform. Th.* **13** (1968) 21–27.
 30 11. N. Cristianini and J. Shawe-Taylor, *An Introduction to Support Vector Machines*
 31 (Cambridge University Press, Cambridge, 2000).
 32 12. C. de Boor, *A Practical Guide to Splines* (Springer-Verlag, New York, 1978).
 33 13. P. A. Devijver and J. Kittler, *Pattern Recognition: a Statistical Approach* (Prentice
 34 Hall, London, 1982).
 35 14. R. Duda, P. Hart and D. G. Stork, *Pattern Classification* (Wiley, New York, 2000).
 36 15. R. A. Fisher, The use of multiple measurements in taxonomic problems. *Ann. Eugen.*
 37 **7** (1936) 179–188.
 38 16. J. Friedman and W. Stuetzle, Projection pursuit regression, *J. Amer. Stat. Assoc.* **76**
 39 (1981) 817–823.
 40 17. J. H. Friedman, On bias, variance, 0-1 loss, and the curse of dimensionality, *Data
 41 Min. Knowl. Discov.* **1** (1997) 55–77.
 42 18. E. Fix and J. L. Hodges, Jr., Discriminatory analysis, nonparametric discrimination,
 43 consistency properties. Randolph Field, Texas, Project 21-49-004, Report No. 4, 1951.
 44 19. K. Fukunaga, *Introduction to Statistical Pattern Recognition* (Academic Press, New
 45 York, 1990).
 46 20. A. K. Ghosh and S. Bose, Backfitting neural networks, *Comput. Stat.* **19** (2004) 193–
 47 210.

- 1 21. A. K. Ghosh and P. Chaudhuri, On data depth and distribution free discriminant
analysis using separating surfaces, *Bernoulli* **11** (2005) 1–27.
- 3 22. D. J. Hand, *Kernel Discriminant Analysis* (Wiley, Chichester, 1982).
- 5 23. T. Hastie, R. Tibshirani and A. Buja, Flexible discriminant analysis, *J. Amer. Stat.*
Assoc. **89** (1994) 1255–1270.
- 7 24. T. Hastie, R. Tibshirani and J. H. Friedman, *The Elements of Statistical Learning:*
Data Mining, Inference and Prediction (Springer Verlag, New York, 2001).
- 9 25. K. Hornik, M. Stinchcombe and H. White, Muli-layer feed forward networks are
universal approximators, *Neural Networks* **2** (1989) 359–366.
- 11 26. P. J. Huber, Projection pursuit (with discussion), *Ann. Stat.* **13** (1985) 435–525.
- 13 27. C. Kooperberg, S. Bose and C. J. Stone, Polychotomus regression, *J. Amer. Stat.*
Assoc. **92** (1997) 117–127.
- 15 28. P. J. M. Laarhoven and E. H. L. Aarts, *Simulated Annealing: Theory and Application*
(D. Reidel Pub, Dordrecht, 1987).
- 17 29. D. J. C. Mackay, A Bayesian framework for backpropagation networks, *Neural Com-*
putation **4** (1992) 448–472.
- 19 30. P. C. Mahalanobis, On the generalized distance in statistics, *Proc. Nat. Acad. Sci.*
India **12** (1936) 49–55.
- 21 31. G. J. McLachlan, *Discriminant Analysis and Statistical Pattern Recognition* (Wiley,
New York, 1992).
- 23 32. S. Mika, G. Ratch, J. Weston, B. Scholkopf, A. J. Smola and K.-R. Muller, *Invariant*
feature Extraction and Classification in Kernel Spaces, Advances in Neural Informa-
tion Processing System, Vol. 1, eds. A. J. Smola, T. K. Leen and K.-R. Muller (MIT
Press, 2000), pp. 526–532.
- 25 33. G. E. Peterson and H. L. Barney, Control methods used in a study of vowels, *J. Acoust.*
Soc. Amer. **24** (1952) 175–185.
- 27 34. B. D. Ripley, Neural networks and related methods for classification (with discussion),
J. Roy. Stat. Soc. Series B **56** (1994) 409–456.
- 29 35. B. D. Ripley, *Pattern Recognition and Neural Networks* (Cambridge University Press,
Cambridge, 1996).
- 31 36. C. B. Roosen and T. Hastie, Automatic smoothing spline projection pursuit, *J. Com-*
put. Graph. Stat. **3** (1994) 235–248.
- 33 37. S. Scholkopf, C. J. C. Burges and A. J. Smola, *Advances in Kernel Methods: Support*
Vector Learning (MIT Press, Cambridge, 1999).
- 35 38. D. W. Scott, *Multivariate Density Estimation: Theory, Practice and Visualization*
(Wiley, New York, 1992).
- 37 39. G. A. F. Seber and C. J. Wild, *Nonlinear Regression* (Wiley, New York, 1989).
- 39 40. B. W. Silverman, *Density Estimation for Statistics and Data Analysis* (Chapman and
Hall, London, 1986).
- 41 41. M. Stone, Cross validation: A review, *Mathematische Operationsforschung und Statis-*
tik, Series Statistics **9** (1977) 127–139.
- 43 42. V. Vapnik, *Statistical Learning Theory* (Wiley, New York, 1998).
- 45 43. M. Zhu, Feature extraction and dimension reduction with applications to classifica-
tion and the analysis of co-occurrence data. Ph.D. Dissertation, Stanford University,
(2001).



Anil Kumar Ghosh received the B.Sc (Hons) degree in statistics from the University of Calcutta in 1996. He received the MStat and PhD degrees from the Indian Statistical Institute, Kolkata, in 1998 and 2004, respectively.

During 2004–2006, he was a postdoctoral research fellow at the Institute of Statistical Science, Academia Sinica, Taiwan and a research associate at the Mathematical Sciences Institute, Australian National University, Canberra. Currently, he is an Assistant Professor at the Mathematics and Statistics Department in Indian Institute of Technology, Kanpur.

The fields of his research interests include pattern recognition, nonparametric and robust statistics, machine learning and statistical computing.



Smarajit Bose received the B. Stat and M. Stat degrees from the Indian Statistical Institute, Calcutta in 1984 and 1986 respectively, and the Ph. D. degree from the University of California, Berkeley in 1992, all in

Statistics.

Between 1991 and 1992, he was a visiting scientist in IBM Almaden Research Center, San Jose where he worked with the Advance Process Monitoring Group. He visited the Statistics Departments of the University of Washington, Seattle during 1992–93 and the Ohio State University, Columbus during 1993–96 where he worked on classification and clustering problems and their applications in Medical Imaging and Ergonomics. In 1996, he joined the Indian Statistical Institute, Calcutta and now is a Professor in the Applied Statistics Division of the Institute. He has also visited the University of California, Santa Barbara in 2001, and during 2002–03.

His current research interests are Pattern Recognition and its applications in Image Processing and Speaker Identification.