



A kernel-based learning algorithm combining kernel discriminant coordinates and kernel principal components

Karol Deręgowski¹, Mirosław Krzyśko^{1,2}

¹President Stanislaw Wojciechowski Higher Vocational State School in Kalisz, Institute of Management, Nowy Świat 4, 62-800 Kalisz, Poland, e-mail: k.deregowski@pwsz.kalisz.pl

²Adam Mickiewicz University, Faculty of Mathematics and Computer Science, Umultowska 87, 61-614 Poznań, Poland, e-mail: mkrzyisko@amu.edu.pl

SUMMARY

Kernel principal components (KPC) and kernel discriminant coordinates (KDC), which are the extensions of principal components and discriminant coordinates, respectively, from a linear domain to a nonlinear domain via the kernel trick, are two very popular nonlinear feature extraction methods. The kernel discriminant coordinates space has proven to be a very powerful space for pattern recognition. However, further study shows that there are still drawbacks in this method. To improve the performance of pattern recognition, we propose a new learning algorithm combining the advantages of KPC and KDC.

Key words: kernel principal components, kernel discriminant coordinates

1. Introduction

Classical principal component analysis (PCA) (Hotelling, 1933) was introduced as a technique for deriving a reduced set of orthogonal linear projections of a single collection of correlated variables $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$, where the projections are ordered by decreasing variance. Principal component analysis is used, for example, in lossy data compression, pattern recognition, and image analysis. In addition to reducing dimensionality, principal component analysis can be used to discover important features of the data. Discovery in principal component analysis takes the form of graphical displays of the principal component scores. The first few principal component scores can reveal whether most of the data actually live on a linear subspace of \mathbb{R}^p , and can

be used to identify outliers, distributional peculiarities, and clusters of points. The last few principal component scores show those linear projections of $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ that have the smallest variance; any principal component with zero or near-zero variance is virtually constant, and hence can be used to detect collinearity, as well as outliers that affect the perceived dimensionality of the data.

When we have samples originating from c groups, we would often like to present them graphically, to see their configuration or to eliminate outlying observations. However it may be difficult to produce such a presentation even if only three features are observed, and with a higher number of features it becomes impossible. A different method must therefore be sought for presenting multidimensional data originating from multiple groups. To make the task easier, in the first step every p -dimensional observation $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ can be transformed into a one-dimensional observation $u_1 = \mathbf{a}_1^T \mathbf{X} = a_{11}X_1 + a_{12}X_2 + \dots + a_{1p}X_p$, and the resulting one-dimensional observations can be presented graphically as points on a line. In the second step we can define a second linear combination $u_2 = \mathbf{a}_2^T \mathbf{X} = a_{21}X_1 + a_{22}X_2 + \dots + a_{2p}X_p$, not correlated with the first, and present the observations graphically as points on a plane.

Generally, the aim is to construct new uncorrelated variables u_1, u_2, \dots, u_s , $s \leq p$, which will be linear combinations of the original observations X_1, X_2, \dots, X_p and which will discriminate the c groups to a maximum degree; that is to say, in the new system the centres of the c groups will be maximally spaced, and the observations from a given group will be maximally concentrated around its centre. These new variables are called discriminant coordinates (see Seber, 1984, p. 270). They are also sometimes called canonical variates, but this name is misleading, because canonical variables with completely different properties occur in canonical correlation analysis. Another name used is “discriminant functions” – this is inappropriate because discriminant functions are surfaces that separate the c groups from one another. The space of discriminant coordinates is a space which is convenient for the

application of various classification methods (methods of discriminant analysis). In the case $c = 2$ we obtain only one discriminant coordinate, coinciding with the well-known linear discriminant function of Fisher (1936).

The linear projection method can be extremely useful in discovering low-dimensional structure when the data actually lie in a linear (or approximately linear) lower-dimensional subspace (called a manifold) M of the input space \mathbb{R}^p . But what can we do if we know or suspect that the data actually lie on a low-dimensional nonlinear manifold, whose structure and dimensionality are both assumed unknown? We can then construct the nonlinear principal components and nonlinear discriminant coordinates.

Kernel principal components (KPC) and kernel discriminant coordinates (KDC), which are the extensions of principal components and discriminant coordinates, respectively, from a linear domain to a nonlinear domain via the kernel trick, are two very popular nonlinear feature extraction methods. The kernel discriminant coordinates space has proven to be a very powerful space for pattern recognition. However, further study shows that there are still drawbacks in this method. One of the major drawbacks of the kernel discriminant coordinate method is that it will lose the within-class scatter information, as for so-called “small sample size” problems, because all of the optimal discriminant vectors in this case are limited to the null space of the within-class scatter matrix – and this information is also important for pattern recognition. To improve the performance of pattern recognition, we propose another learning algorithm combining the advantages of KPC and KDC. Our proposed algorithm can be divided into three steps:

- (1) compute the optimal discriminant vectors of KDC;
- (2) compute the optimal vectors of KPC;
- (3) use the two kinds of features for recognition.

The paper is organized as follows. In section 2, classical principal components are presented. Kernel principal components are presented in section 3. In section 4, classical discriminant coordinates are presented, and kernel discriminant coordinates are described in section 5. The discriminant

algorithm based on a combination of features from the kernel discriminant coordinates space and kernel principal components space is described in section 6. Finally, section 7 examines the quality of the new discriminant algorithm presented in this paper.

2. Classical principal component analysis

Assume that the random p -vector $\mathbf{X} = (X_1, X_2, \dots, X_p)^T$ has mean $\boldsymbol{\mu}$ and the $(p \times p)$ covariance matrix $\boldsymbol{\Sigma}$. PCA seeks to replace the set of p (unordered and correlated) input variables, X_1, X_2, \dots, X_p by a (potentially smaller) set of t (ordered and uncorrelated) linear projections, ξ_1, \dots, ξ_t ($t \leq p$), of the input variables,

$$\xi_j = \mathbf{b}_j^T \mathbf{X} = b_{j1}X_1 + \dots + b_{jp}X_p, j = 1, 2, \dots, t; \quad (2.1)$$

where we minimize the loss of information due to the replacement.

In PCA, “information” is interpreted as the “total variation” of the original input variables,

$$\sum_{j=1}^p \text{Var}(X_j) = \text{tr}(\boldsymbol{\Sigma}).$$

From the spectral decomposition theorem, we can write $\boldsymbol{\Sigma} = \mathbf{U}\boldsymbol{\Lambda}\mathbf{U}^T$, $\mathbf{U}^T\mathbf{U} = \mathbf{I}_p$, where the diagonal matrix $\boldsymbol{\Lambda}$ has as diagonal elements the eigenvalues $\{\lambda_j\}$ of $\boldsymbol{\Sigma}$, and the columns of \mathbf{U} are the eigenvectors of $\boldsymbol{\Sigma}$. Thus the total variation is $\text{tr}(\boldsymbol{\Sigma}) = \text{tr}(\boldsymbol{\Lambda}) = \sum_{j=1}^p \lambda_j$.

The j th coefficient vector, $\mathbf{b}_j = (b_{j1}, \dots, b_{jp})^T$, is chosen so that:

- The first t linear projections ξ_j , $j = 1, 2, \dots, t$, of \mathbf{X} are ranked in importance according to their variances $\{\text{Var}(\xi_j)\}$, which are listed in decreasing order of magnitude:

$$\text{Var}(\xi_1) \geq \text{Var}(\xi_2) \geq \dots \geq \text{Var}(\xi_t).$$

- ξ_j is uncorrelated with all ξ_k , $k < j$.

The linear projections (1) are then known as the first t principal components of \mathbf{X} .

In practice, we estimate the principal components using N independent observations, $\{\mathbf{X}_i, i = 1, 2, \dots, N\}$, on \mathbf{X} . We estimate $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ by $\hat{\boldsymbol{\mu}} = \bar{\mathbf{X}} = N^{-1} \sum_{i=1}^N \mathbf{X}_i$ and

$$\hat{\boldsymbol{\Sigma}} = \mathbf{S} = (N-1)^{-1} \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T, \quad (2.2)$$

respectively.

The ordered eigenvalues of $\hat{\boldsymbol{\Sigma}}$ are denoted by $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_p \geq 0$, and the eigenvector associated with the j th largest sample eigenvalue $\hat{\lambda}_j$ is the j th sample eigenvector $\hat{\mathbf{v}}_j, j = 1, 2, \dots, p$.

The coordinates of the projection of the i th observation \mathbf{X}_i on the j th principal components are equal to:

$$\hat{\xi}_{ij} = \hat{\mathbf{v}}_j^T \mathbf{X}_i, \quad (2.3)$$

$i = 1, 2, \dots, N, j = 1, 2, \dots, p$.

A sample measure of how well the first t principal components represent the p original variables is given by the statistic

$$\frac{\hat{\lambda}_1 + \dots + \hat{\lambda}_t}{\hat{\lambda}_1 + \dots + \hat{\lambda}_p}$$

which is the proportion of the total sample variation that is explained by the first t sample principal components.

It is hoped that the sample variances of the first few sample PCs will be large, and that the remainder will be small enough for the corresponding set of sample PCs to be omitted. A variable that does not change much (relative to other variables) in independent measurements may be treated approximately as a constant, and so omitting such low-variance sample PCs and focusing exclusively on the high-variance sample PCs is a convenient way of reducing the dimensionality of the data set. For diagnostic and data analytic purposes, it

is usual to plot the first sample PC scores against the second sample PC scores, $(\hat{\xi}_{i1}, \hat{\xi}_{i2})$, where $\hat{\xi}_{ij}$ is given by (2.3), for $i = 1, \dots, N, j = 1, 2$.

3. Kernel principal component analysis

An approach that generalizes linear PCA is that of kernel PCA (Schölkopf et al., 1998). This is an application of what are called kernel methods.

Let $\mathbf{X}_i \in \mathbf{R}^p$, $i = 1, 2, \dots, N$, be the input data points. We can think of kernel PCA as a two-step process:

- Nonlinearly transform the i th input data point $\mathbf{X}_i \in \mathbf{R}^p$ into a point $\Phi(\mathbf{X}_i)$ in an N_H -dimensional feature space H (the Hilbert space), where

$$\Phi(\mathbf{X}_i) = (\Phi_1(\mathbf{X}_i), \dots, \Phi_{N_H}(\mathbf{X}_i))^T \in H, i = 1, 2, \dots, N.$$

- The transformation $\Phi : \mathbf{R}^p \rightarrow H$ is called a feature transformation, and each $\{\Phi_j\}$ is a nonlinear transformation.
- Given $\Phi(\mathbf{X}_1), \dots, \Phi(\mathbf{X}_N) \in H$, solve a linear PCA problem in the feature space H , which has a higher dimensionality than that of the input space (i.e. $N_H > p$).

Consider the data presented in Figure 1. Let $\mathbf{X}_i = (X_{i1}, X_{i2})^T$, and define $\Phi : \mathbf{R}^2 \rightarrow \mathbf{R}^3$ by $\Phi(\mathbf{X}_i) = \Phi(X_{i1}, X_{i2}) = (X_{i1}^2, \sqrt{2}X_{i1}X_{i2}, X_{i2}^2)^T = (z_{i1}, z_{i2}, z_{i3})^T$.

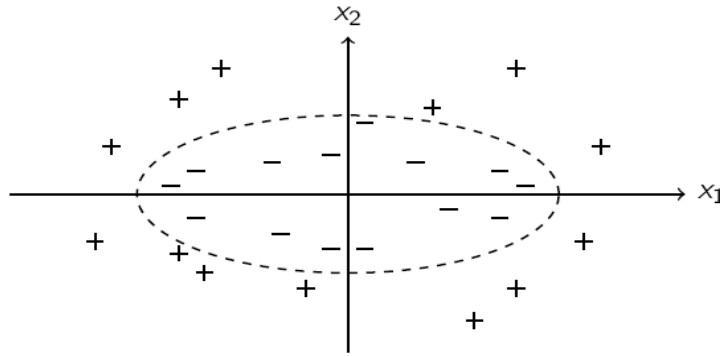


Figure 1. Original data in the plane. The data cannot be separated linearly

With this Φ , a difficult nonlinear classification problem in \mathbb{R}^2 is converted to a standard linear classification task in \mathbb{R}^3 (see Figure 2).

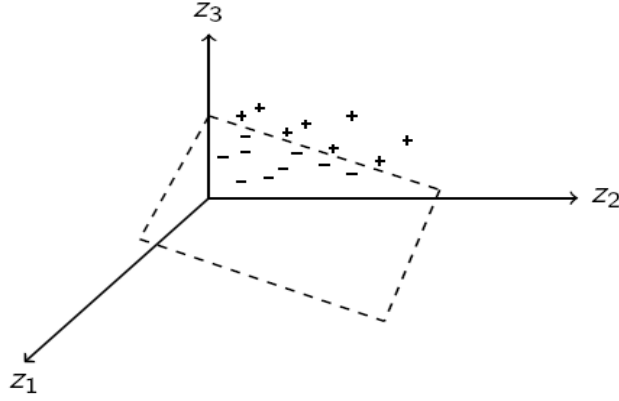


Figure 2. A three-dimensional representation of the pluses and minuses

Let $\mathbf{X}_i = (X_{i1}, X_{i2})^T$ and $\mathbf{Y}_i = (Y_{i1}, Y_{i2})^T$ be two vectors in the input space \mathbb{R}^2 , and consider the transformation to \mathbb{R}^3 used earlier. Let $\Phi(\mathbf{X}_i)$ and $\Phi(\mathbf{Y}_i)$ be two feature vectors generated by \mathbf{X}_i and \mathbf{Y}_i . Now consider the inner product $\Phi^T(\mathbf{X}_i) \Phi(\mathbf{Y}_i)$ in the feature space. It takes the form

$$\begin{aligned} \Phi^T(\mathbf{X}_i) \Phi(\mathbf{Y}_i) &= (X_{i1}^2, \sqrt{2}X_{i1}X_{i2}, X_{i2}^2)(Y_{i1}^2, \sqrt{2}Y_{i1}Y_{i2}, Y_{i2}^2)^T = \\ &= (X_{i1}Y_{i1} + X_{i2}Y_{i2})^2 = (\mathbf{X}_i^T \mathbf{Y}_i)^2 = k(\mathbf{X}_i, \mathbf{Y}_i). \end{aligned} \quad (3.1)$$

Equation (4) shows how an inner product based on Φ converts to a function of the two inputs. Since choosing an inner product and performing computations with it in the feature space can quickly become computationally infeasible, it would be desirable to choose a function k , called a kernel, so as to summarize the geometry of the feature space vectors and ignore Φ entirely.

Now the kernel trick can be applied. Suppose a function $k(\cdot, \cdot) : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ operating on the input space can be found such that the feature space inner products are computed directly through k as in (4). Then explicit use of Φ has been avoided, and yet results can be obtained as if Φ were used. This direct

computation of feature space inner products without explicitly manipulating the feature space vectors themselves is known as the kernel trick.

The existence of the transformation $\Phi : \mathbb{R}^p \rightarrow H$ such that

$$\Phi^T(\mathbf{X}_i)\Phi(\mathbf{Y}_i) = k(\mathbf{X}_i, \mathbf{Y}_i)$$

is guaranteed by the following theorem.

Theorem 3.1 (Moore–Aronszajn) (Aronszajn, 1950). Let $k : \mathbb{R}^p \times \mathbb{R}^p \rightarrow \mathbb{R}$ be a bivariate symmetric continuous real-valued function and H_k be a reproducing kernel Hilbert space (RKHS). Then there exists a transformation $\Phi : \mathbb{R}^p \rightarrow H_k$ such that

$$k(\mathbf{X}_i, \mathbf{Y}_i) = \Phi^T(\mathbf{X}_i)\Phi(\mathbf{Y}_i)$$

if and only if the matrix $\mathbf{K} = (k_{ij})$ is nonnegative definite, where $k_{ij} = k(\mathbf{X}_i, \mathbf{X}_j)$, $i, j = 1, \dots, N$.

The matrix \mathbf{K} is known as the kernel matrix. For a given bivariate function k , verifying the above conditions might not be easy. In practice, there exist many functions that have been shown to be valid kernels, and fortunately many of them deliver good performance on real-world data.

A short annotated list is presented in Table 1.

Table 1. Kernel functions

Kernel	$k(\mathbf{x}, \mathbf{y})$
Homogeneous polynomial kernel	$(\mathbf{x}^T \mathbf{y})^d$, d is an integer
Inhomogeneous polynomial kernel	$(\mathbf{x}^T \mathbf{y} + c)^d$, $c > 0$
Gaussian radial basis function	$\exp(-c \ \mathbf{x} - \mathbf{y}\ ^2)$, $c > 0$
Laplacian	$\exp(-c \ \mathbf{x} - \mathbf{y}\)$, $c > 0$

In order to carry out linear PCA in the feature space so that it mimics the standard treatment of PCA (as carried out in the input space), we have to find eigenvalues $\gamma \geq 0$ and nonzero eigenvectors $\mathbf{u} \in H$ of the estimated covariance matrix

$$\mathbf{C} = \frac{1}{N-1} \sum_{i=1}^N \Phi(\mathbf{X}_i) \Phi^T(\mathbf{X}_i) \quad (3.2)$$

of the centred and nonlinearly transformed input vectors. The eigenequation $\mathbf{C}\mathbf{u} = \gamma\mathbf{u}$, where \mathbf{u} is the eigenvector corresponding to the eigenvalue $\gamma \geq 0$ of \mathbf{C} , can be written in the equivalent form

$$\Phi^T(\mathbf{X}_i) \mathbf{C}\mathbf{u} = \gamma \Phi^T(\mathbf{X}_i) \mathbf{u}, \quad i = 1, 2, \dots, N. \quad (3.3)$$

Because

$$\mathbf{C}\mathbf{u} = \frac{1}{N-1} \sum_{i=1}^N \Phi(\mathbf{X}_i) \Phi^T(\mathbf{X}_i) \mathbf{u}$$

all solutions \mathbf{u} with nonzero eigenvalue γ are contained in the span of $\Phi(\mathbf{X}_1), \dots, \Phi(\mathbf{X}_N)$. Hence there exist coefficients α_k , $k = 1, 2, \dots, N$, such that

$$\mathbf{u} = \sum_{k=1}^N \alpha_k \Phi(\mathbf{X}_k). \quad (3.4)$$

Substituting (3.4) for \mathbf{u} in (3.3), we obtain that

$$\frac{1}{N-1} \sum_{j=1}^N \alpha_j \Phi^T(\mathbf{X}_i) \sum_{k=1}^N \Phi(\mathbf{X}_k) \Phi^T(\mathbf{X}_k) \Phi(\mathbf{X}_j) = \gamma \sum_{k=1}^N \alpha_k \Phi^T(\mathbf{X}_i) \Phi(\mathbf{X}_k), \quad (3.5)$$

for all $i = 1, 2, \dots, N$.

The eigenequation (3.5) can be written as $\mathbf{K}^2 \boldsymbol{\alpha} = N \gamma \mathbf{K} \boldsymbol{\alpha}$, where $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_N)^T$, or as

$$\mathbf{K}^2 \boldsymbol{\alpha} = \tilde{\gamma} \mathbf{K} \boldsymbol{\alpha}, \quad (3.6)$$

where $\tilde{\gamma} = (N-1) \gamma$, $\mathbf{K} = (k_{ij})$ and $k_{ij} = k(\mathbf{X}_i, \mathbf{X}_j) = \Phi^T(\mathbf{X}_i) \Phi(\mathbf{X}_j)$, $i, j = 1, 2, \dots, N$.

To find solutions of (3.6), we solve the eigenvalue problem

$$\mathbf{K}\boldsymbol{\alpha} = \tilde{\gamma}\boldsymbol{\alpha} \quad (3.7)$$

for nonzero eigenvalues. Clearly, all solutions of (3.7) also satisfy (3.6). Moreover, it can be shown that any additional solutions of (3.7) do not make a difference in the expansion (3.4) and thus are not of interest to us.

We assumed earlier that the vectors $\{\Phi(\mathbf{X}_i)\}$, $i = 1, 2, \dots, N$, are centred. In the general case we cannot centre the vectors $\{\Phi(\mathbf{X}_i)\}$, because we do not know the form of the function Φ . Let

$$\tilde{\Phi}(\mathbf{X}_i) = \Phi(\mathbf{X}_i) - \frac{1}{N} \sum_{k=1}^N \Phi(\mathbf{X}_k)$$

and

$$\tilde{\mathbf{K}} = (\tilde{k}(\mathbf{X}_i, \mathbf{X}_j)) = (\langle \tilde{\Phi}(\mathbf{X}_i), \tilde{\Phi}(\mathbf{X}_j) \rangle)$$

$i, j = 1, 2, \dots, N$. We cannot compute the matrix $\tilde{\mathbf{K}}$ directly, but we can express it in terms of the matrix \mathbf{K} : $\tilde{\mathbf{K}} = \mathbf{P}\mathbf{K}\mathbf{P}$, where $\mathbf{P} = \mathbf{I}_N - N^{-1}\mathbf{1}_N\mathbf{1}_N^T$ is the centring matrix.

Hence in the general case the construction of kernel principal components must be based on the matrix $\tilde{\mathbf{K}}$. Because the kernel matrix \mathbf{K} is nonnegative definite, $\tilde{\mathbf{K}}$ is nonnegative definite also. This results from the fact (Seber, 1984, p. 521) that if $\mathbf{A} \geq 0$, then $\mathbf{C}\mathbf{A}\mathbf{C}^T \geq 0$. In our case $\mathbf{C} = \mathbf{P}$, where \mathbf{P} is a symmetric matrix, i.e. $\mathbf{P} = \mathbf{P}^T$. Hence, all eigenvalues of $\tilde{\mathbf{K}}$ are nonnegative.

$$\text{The coordinates of the projection of the data matrix } \begin{bmatrix} \Phi^T(\mathbf{X}_1) \\ \vdots \\ \Phi^T(\mathbf{X}_N) \end{bmatrix}$$

on the j th kernel principal components are equal to $\tilde{\mathbf{K}}\hat{\mathbf{a}}_j$, $j = 1, 2, \dots, N$.

4. Classical discriminant coordinates

Consider a set of N p -dimensional observations $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$, where the vectors $\mathbf{X}_i \in \mathbb{R}^p$ are grouped in c disjoint classes, and every \mathbf{X}_i belongs to one and only one class.

Let $V = \{1, 2, \dots, N\}$ be the set of indices of the set of observations $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N$, and let V be partitioned into c disjoint subsets V_i such that

$V_i \cap V_j = \emptyset$ for $i \neq j$, $\bigcup_{j=1}^c V_j = V$, and each V_j contains n_j elements such that $\sum_{j=1}^c n_j = N$.

Let $\bar{\mathbf{X}} = N^{-1} \sum_{i=1}^N \mathbf{X}_i$ be the sample mean, and let $\bar{\mathbf{X}}_j = n_j^{-1} \sum_{i \in V_j} \mathbf{X}_i$ be the mean of the j th class, for $j=1, 2, \dots, c$.

We denote by \mathbf{S}_t the total scatter matrix, and by \mathbf{S}_b the between-class scatter matrix. We have

$$\mathbf{S}_t = \sum_{i=1}^N (\mathbf{X}_i - \bar{\mathbf{X}})(\mathbf{X}_i - \bar{\mathbf{X}})^T,$$

$$\mathbf{S}_b = \sum_{j=1}^c n_j (\bar{\mathbf{X}}_j - \bar{\mathbf{X}})(\bar{\mathbf{X}}_j - \bar{\mathbf{X}})^T.$$

We seek a set of vectors $\mathbf{a}_i \in \mathbb{R}^p$ which maximize the measure of separation of the c classes

$$J(\mathbf{a}) = \frac{\mathbf{a}^T \mathbf{S}_b \mathbf{a}}{\mathbf{a}^T \mathbf{S}_t \mathbf{a}},$$

subject to the additional restriction $\mathbf{a}_i^T \mathbf{S}_t \mathbf{a}_j = \delta_{ij}$ (Kronecker delta), which means that the variables $u_i = \mathbf{a}_i^T \mathbf{X}$, called discriminant coordinates, are uncorrelated. The \mathbf{a}_i are called directional vectors. Finding the directional vectors reduces to solving the following generalized eigenvalue problem:

$$\mathbf{S}_b \mathbf{a}_i = \lambda_i \mathbf{S}_t \mathbf{a}_i, \quad i=1, 2, \dots, s, \quad (4.1)$$

where $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s \geq \lambda_{s+1} = 0$ and $s = \text{rank}(\mathbf{S}_b) \leq \min\{p, c-1\}$.

Because $\mathbf{S}_t = \mathbf{S}_b + \mathbf{S}_w$, where \mathbf{S}_w is the pooled within-class scatter matrix, finding classical discriminant coordinates is equivalent to finding a solution to the generalized eigenvalue problem in the form $\mathbf{S}_b \mathbf{a} = \lambda/(1 - \lambda) \mathbf{S}_w \mathbf{a}$.

Using matrix notation, the generalized problem (4.1) can be given in the form

$$\mathbf{S}_b \mathbf{A} = \mathbf{S}_t \mathbf{A} \mathbf{\Lambda}, \quad (4.2)$$

where $\mathbf{A} = (\mathbf{a}_1, \dots, \mathbf{a}_s)$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_s)$.

If \mathbf{S}_t is a non-singular matrix, we obtain $\mathbf{S}_t^{-1}\mathbf{S}_b\mathbf{A} = \mathbf{A}\mathbf{\Lambda}$, which is an ordinary eigenvalue problem.

Hence the pair $(\lambda_i, \mathbf{a}_i)$, where λ_i is an eigenvalue of the matrix $\mathbf{S}_t^{-1}\mathbf{S}_b$ and \mathbf{a}_i is the corresponding eigenvector, can be used to construct discriminant coordinates. Because $\text{rank}(\mathbf{S}_b) \leq c-1$, we obtain a discriminant coordinate space with dimension at least $c-1$.

5. Kernel discriminant coordinates

Kernel discriminant coordinates were introduced independently by Mika et al. (1999) and by Baudat and Anouar (2000). The method is described in a book by Shawe-Taylor and Cristianini (2004).

The original space \mathbf{R}^p is transformed non-linearly into a feature space H_k

$$\Phi : \mathbf{R}^p \rightarrow H_k,$$

where Φ is a vector-valued function, and H_k is a reproducing kernel Hilbert space (RKHS).

The vector $\Phi(\mathbf{X}_i) = \tilde{\mathbf{X}}_i$ is called the feature vector corresponding to the observation $\mathbf{X}_i \in \mathbf{R}^p$, $i = 1, 2, \dots, N$. The non-linear transformation Φ is in general not known; however we select a known form of the non-negative definite kernel function

$$k(\mathbf{X}, \mathbf{Y}) = \Phi^T(\mathbf{X})\Phi(\mathbf{Y}) = \tilde{\mathbf{X}}^T\tilde{\mathbf{Y}}.$$

Let $\tilde{\mathbf{S}}_b$ and $\tilde{\mathbf{S}}_t$ denote respectively the between-class and total scatter matrices in the feature space. We have

$$\tilde{\mathbf{S}}_t = \sum_{i=1}^N (\tilde{\mathbf{X}}_i - \tilde{\mathbf{X}})(\tilde{\mathbf{X}}_i - \tilde{\mathbf{X}})^T,$$

$$\tilde{\mathbf{S}}_b = \sum_{j=1}^c n_j (\tilde{\mathbf{X}}_j - \tilde{\mathbf{X}})(\tilde{\mathbf{X}}_j - \tilde{\mathbf{X}})^T,$$

where $\tilde{\mathbf{X}}_k$ and $\tilde{\mathbf{X}}$ are respectively the mean of the k th class and total mean in the feature space. Finding discriminant coordinates in the feature space reduces to solving the following optimization problem:

$$\mathbf{b}_i = \arg \max \frac{\mathbf{b}^T \tilde{\mathbf{S}}_b \mathbf{b}}{\mathbf{b}^T \tilde{\mathbf{S}}_t \mathbf{b}}, \quad \mathbf{b}_i^T \tilde{\mathbf{S}}_b \mathbf{b}_j = \delta_{ij}. \quad (5.3)$$

We know that there exist coefficients b_{ij} such that $\mathbf{b}_i = \sum_{j=1}^n b_{ij} \tilde{\mathbf{X}}_j$. Hence the optimization problem (5.3) is equivalent to the problem

$$\mathbf{b}_i = \arg \max \frac{\mathbf{b}^T \tilde{\mathbf{K}} \mathbf{D} \tilde{\mathbf{K}} \mathbf{b}}{\mathbf{b}^T \tilde{\mathbf{K}} \tilde{\mathbf{K}} \mathbf{b}}, \quad \mathbf{b}_i^T \tilde{\mathbf{K}} \tilde{\mathbf{K}} \mathbf{b}_j = \delta_{ij}$$

and the optimum vectors $\mathbf{b}_1, \dots, \mathbf{b}_s$ are equal to the eigenvectors corresponding to the maximum eigenvalues in the generalized problem

$$\tilde{\mathbf{K}} \mathbf{D} \tilde{\mathbf{K}} \mathbf{b} = \lambda \tilde{\mathbf{K}} \tilde{\mathbf{K}} \mathbf{b}, \quad \tilde{\mathbf{K}} = \mathbf{P} \mathbf{K} \mathbf{P}, \quad (5.4)$$

where $\mathbf{K} = (k_{ij})$ ($k_{ij} = k(\mathbf{X}_i, \mathbf{X}_j)$) is the kernel matrix, $\mathbf{P} = \mathbf{I}_N - N^{-1} \mathbf{1}_N \mathbf{1}_N^T$ is the centring matrix, and the matrix \mathbf{D} is defined by

$$D_{ij} = \begin{cases} \frac{1}{n_k}, & \text{if } \mathbf{X}_i \text{ and } \mathbf{X}_j \text{ belong to the } k\text{th class,} \\ 0, & \text{otherwise.} \end{cases}$$

Let $\mathbf{B} = [\mathbf{b}_1, \dots, \mathbf{b}_s]$ and $\mathbf{\Lambda} = \text{diag}(\lambda_1, \dots, \lambda_s)$. In matrix notation the problem of (5.4) has the form

$$\tilde{\mathbf{K}} \mathbf{D} \tilde{\mathbf{K}} \mathbf{B} = \tilde{\mathbf{K}} \tilde{\mathbf{K}} \mathbf{B} \mathbf{\Lambda}. \quad (5.5)$$

Solving the generalized eigenvalue problem presents certain difficulties, because both matrices are non-negative definite.

One of the ways of obtaining an approximate solution to (5.5), also used in ridge regression, is to regularize the matrix $\tilde{\mathbf{K}} \tilde{\mathbf{K}}$, that is to replace $\tilde{\mathbf{K}} \tilde{\mathbf{K}}$ with a new non-singular matrix $\tilde{\mathbf{K}} \tilde{\mathbf{K}} + \varepsilon \mathbf{I}$ (Friedman, 1989; Mika et al., 1999). We then solve the new generalized eigenvalue problem

$$\tilde{\mathbf{K}} \mathbf{D} \tilde{\mathbf{K}} \mathbf{B} = (\tilde{\mathbf{K}} \tilde{\mathbf{K}} + \varepsilon \mathbf{I}) \mathbf{B} \mathbf{\Lambda}.$$

Because the matrix $\tilde{\mathbf{K}}\tilde{\mathbf{K}} + \varepsilon\mathbf{I}$ is non-singular, the above problem can be reduced to the classical eigenvalue problem

$$(\tilde{\mathbf{K}}\tilde{\mathbf{K}} + \varepsilon\mathbf{I})^{-1}\tilde{\mathbf{K}}\mathbf{D}\tilde{\mathbf{K}}\mathbf{B} = \mathbf{B}\boldsymbol{\Lambda},$$

whose solution has the form

$$\mathbf{B} = \text{eig}[(\tilde{\mathbf{K}}\tilde{\mathbf{K}} + \varepsilon\mathbf{I})^{-1}\tilde{\mathbf{K}}\mathbf{D}\tilde{\mathbf{K}}].$$

Hence the pair $(\lambda_i, \mathbf{b}_i)$, where λ_i is an eigenvalue of the matrix $(\tilde{\mathbf{K}}\tilde{\mathbf{K}} + \varepsilon\mathbf{I})^{-1}\tilde{\mathbf{K}}\mathbf{D}\tilde{\mathbf{K}}$ and \mathbf{b}_i is the corresponding eigenvector, can be used to construct kernel discriminant coordinates. Because $\text{rank}(\tilde{\mathbf{K}}\mathbf{D}\tilde{\mathbf{K}}) \leq c-1$, we obtain a discriminant coordinate space with dimension at least $c-1$. The coordinates of the projection of the data matrix

$$\begin{bmatrix} \tilde{\boldsymbol{\Phi}}^T(\mathbf{X}_1) \\ \vdots \\ \tilde{\boldsymbol{\Phi}}^T(\mathbf{X}_N) \end{bmatrix}$$

on the j th kernel principal components are equal to $\tilde{\mathbf{K}}\hat{\boldsymbol{\alpha}}_j, j = 1, 2, \dots, c-1$.

6. A new discriminant space

Consider the $(c-1)$ -dimensional space of the first kernel discriminant variables. The training sample $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$ transformed into this space will be denoted by $\{\mathbf{Y}_1, \mathbf{Y}_2, \dots, \mathbf{Y}_N\}$, where $\mathbf{Y}_i \in \mathbb{R}^{c-1}$. We then further consider the $(c-1)$ -dimensional space of the first kernel principal components. The training sample transformed into this space will be denoted by $\{\mathbf{Z}_1, \mathbf{Z}_2, \dots, \mathbf{Z}_N\}$, where $\mathbf{Z}_i \in \mathbb{R}^{c-1}$.

We will create a new $2(c-1)$ -dimensional space determined by the first $(c-1)$ kernel discriminant variables and the first $(c-1)$ kernel principal components. The directional vectors determining this new space are normed so that their length is equal to 1. The new space combines the advantages of the space of kernel discriminant coordinates and the space of kernel principal components. In this space it is possible to apply a variety of classification algorithms, obtaining an improvement in classification quality.

7. Example

Three types of discriminant spaces were considered: the $(c-1)$ -dimensional space of kernel principal components described in section 3, the $(c-1)$ -dimensional space of kernel discriminant coordinates described in section 5, and the combined $2(c-1)$ -dimensional space of kernel discriminant coordinates and kernel principal components.

The kernel spaces were built using the Gaussian kernel

$$k(\mathbf{x}, \mathbf{y}) = \exp(-d \|\mathbf{x} - \mathbf{y}\|^2),$$

where d is a positive constant. We consider the lower triangle in the table of squared distances between elements of the training sample $\{\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_N\}$. The value of d was taken to be the reciprocal of the arithmetic mean of the elements of the lower triangle. The matrix $\tilde{\mathbf{K}}\tilde{\mathbf{K}} + \varepsilon\mathbf{I}$ was taken with $\varepsilon = 10^{-5}$.

To check the usefulness of the three aforementioned discriminant spaces in the process of classification, each space was used to classify objects belonging to 25 different training samples from the University of California database (Bache and Lichman, 2013).

Classification was performed by the method of linear discriminant functions. Table 2 shows that the largest percentage of correct classifications is achieved in the combined space of kernel discriminant coordinates and kernel principal components. This is the case for 23 of the training samples. For one training sample (Ring) we have the same percentage of correct classifications in the kernel discriminant coordinate space as in the combined space of kernel discriminant coordinates and kernel principal components. For one training sample (Hayes-Roth) the kernel discriminant coordinate space gives better results than the combined space. The results imply that the combined space should be chosen as the best discriminant space.

Table 2. Comparison of percentages of correct classification

Data	N	c	p	Percentage of correct classification		
Balance	625	3	4	87.36 (546/625)	89.12 (557/625)	90.88 (568/625)
Car	1728	4	6	69.56 (1202/1728)	83.33 (1440/1728)	86.46 (1494/1728)
Choice	1473	3	9	45.76 (674/1473)	45.96 (677/1473)	54.11 (797/1473)
Danesym	300	3	2	67.33 (202/300)	68.33 (205/300)	71.00 (213/300)
Danford	45	4	3	42.22 (19/45)	44.44 (20/45)	53.33 (24/45)
Dermatology	358	6	34	75.14 (269/358)	83.80 (300/358)	95.25 (341/358)
EnVowel	990	11	10	75.96 (752/990)	93.84 (929/990)	95.86 (949/990)
Glass	214	6	9	62.15 (133/214)	64.49 (138/214)	73.83 (158/214)
HayesRoth	132	3	3	56.82 (75/132)	85.61 (113/132)	80.30 (106/132)
Image	2310	7	19	76.84 (1637/2310)	78.40 (1811/2310)	87.06 (2011/2310)
Lenses	24	3	4	70.83 (17/24)	62.50 (15/24)	75.00 (18/24)
Oil	56	3	5	76.79 (40/56)	76.79 (40/56)	80.36 (45/56)
PIVowel	1130	6	5	62.74 (709/1130)	63.98 (723/1130)	73.45 (830/1130)
Ring	1000	2	2	65,90 (659/1000)	71,70 (717/1000)	71,70 (717/1000)
Risk	87	3	2	97.70 (85/87)	97.70 (85/87)	100.00 (87/87)
Smith	45	4	3	55.56 (25/45)	60.00 (27/45)	66.67 (30/45)
Soybean	47	4	35	85.11(40/47)	87.23(41/47)	95.74(45/47)
Tae	151	3	5	41.72 (63/151)	45.70 (69/151)	46.36 (70/151)
Three-class	300	3	2	90.00 (270/300)	91.33 (274/300)	92.67 (278/300)
Thyroid	215	3	5	86.51 (186/215)	85.58 (184/215)	93.49 (201/215)
Wale-form	900	3	21	85.11 (766/900)	88.33 (795/900)	93.67 (843/900)
Wine	178	3	13	70.22 (125/178)	71.34 (127/178)	71.91 (128/178)
WineQuality	1599	6	11	49.22 (787/1599)	49.53 (792/1599)	51.34 (821/1599)
Vehicle	846	4	18	43.38 (367/846)	47.28 (400/846)	47.75 (404/846)
Zoo	101	7	16	96.04 (97/101)	96.04 (97/101)	99.01 (100/101)

REFERENCES

- Aronszajn N. (1950): Theory of reproducing kernels. *Transactions of the American Mathematical Society* 68: 337–404.
- Badat G., Anouar F. (2000): Generalized discriminant analysis using a kernel approach. *Neural Computation* 12: 2385–2404.
- Bache K., Lichman M. (2013): UCI Machine Learning Repository [<http://archive.ics.uci.edu/ml>]. Irvine, CA: University of California, School of Information and Computer Science.

- Fisher R.A. (1936): The use of multiple measurements in taxonomic problem. *Annals of Eugenics* 7: 179–188.
- Friedman J.H. (1989): Regularized discriminant analysis. *Journal of the American Statistical Association* 84: 165–175.
- Hotelling H. (1933): Analysis of a complex of statistical variables into principal components. *Journal of Educational Psychology* 24: 417–441, 498–520.
- Mika S., Rätsch G., Weston J., Schölkopf B., Müller K.R. (1999): Fisher discriminant analysis with kernels. In Y.H. Hu, J. Larsen, E. Wilson, and S. Douglas (eds.), *Neural Networks for Signal Processing IV*: 41–48.
- Schölkopf B., Smola A., Müller K.B. (1998): Nonlinear component analysis as a kernel eigenvalues problem. *Neural Computation* 10: 1299–1319.
- Seber G.A.F. (1984): *Multivariate Observations*. Wiley, New York.
- Shawe-Taylor J., Cristianini N. (2004): *Kernel methods for pattern analysis*. Cambridge University Press, Cambridge, UK.