

TILM3592
Advanced Statistical Learning
Kernelized Linear Discriminant Analysis

NAME
MATR.

1 Introduction

This paper will try to explain the generalization of linear discriminant analysis to a nonlinear, kernelized version. Focus is perhaps not put as much on the classification as on the dimension-reduction part, focusing perhaps more on kernelizing Fisher's linear discriminant analysis so that the resulting projection can be used as a basis for a classifier. It was partly inspired by exercise 12.10, *Kernels and linear discriminant analysis*, in [2].

2 (Fisher) Linear Discriminant Analysis

Consider trying to classify two classes by separating them using a hyperplane, similar to the (linear, soft or hard margin) support vector machine or the logistic regression algorithm. Another approach could be to represent the two classes using their means and covariances. The problem could then take the form of finding a projection of the data such that the spread between the two classes (means) is maximized while the spread within the groups is minimized. This can be likened to principal components analysis where one instead tries to find a projection that maximizes the spread and we will see that this approach is similar in some ways.

To make the previous more concrete, let \mathbf{X} be a $(p \times 1)$ random vector and \mathbf{w} be a $(p \times 1)$ vector representing a linear combination, then consider the covariance and mean of the random variable $\mathbf{w}^\top \mathbf{X}$, $\text{Cov}(\mathbf{w}^\top \mathbf{X}, \mathbf{w}^\top \mathbf{X})$ and $\mathbb{E}(\mathbf{w}^\top \mathbf{X})$. Since the expectation operator is linear we know that $\mathbb{E}(\mathbf{w}^\top \mathbf{X}) = \mathbf{w}^\top \mathbb{E}(\mathbf{X})$ and since

$$\text{Cov}(\mathbf{X}, \mathbf{X}) := \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] = \Sigma_{\mathbf{X}\mathbf{X}},$$

we get

$$\begin{aligned}
\underline{\text{Cov}(\mathbf{w}^\top \mathbf{X}, \mathbf{w}^\top \mathbf{X})} &= \mathbb{E}[(\mathbf{w}^\top \mathbf{X} - \mathbb{E}[\mathbf{w}^\top \mathbf{X}])(\mathbf{w}^\top \mathbf{X} - \mathbb{E}[\mathbf{w}^\top \mathbf{X}])^\top] \\
&= \mathbb{E}[(\mathbf{w}^\top \mathbf{X} - \mathbf{w}^\top \mathbb{E}[\mathbf{X}])(\mathbf{w}^\top \mathbf{X} - \mathbf{w}^\top \mathbb{E}[\mathbf{X}])^\top] \\
&= \mathbb{E}[\mathbf{w}^\top (\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top \mathbf{w}] \\
&= \mathbf{w}^\top \mathbb{E}[(\mathbf{X} - \mathbb{E}[\mathbf{X}])(\mathbf{X} - \mathbb{E}[\mathbf{X}])^\top] \mathbf{w} = \underline{\mathbf{w}^\top \Sigma_{\mathbf{X}\mathbf{X}} \mathbf{w}}.
\end{aligned}$$

Just as a sanity check, since \mathbf{X} is $(p \times 1)$ we know that $\text{Cov}(\mathbf{X}, \mathbf{X})$ is a $(p \times p)$ matrix and since \mathbf{w}^\top is $(1 \times p)$ the random variable $\mathbf{w}^\top \mathbf{X}$ is (1×1) (scalar) the covariance (variance) $\text{Cov}(\mathbf{w}^\top \mathbf{X}, \mathbf{w}^\top \mathbf{X})$ also has to be scalar. This is true since $\mathbf{w}^\top \Sigma_{\mathbf{X}\mathbf{X}} \mathbf{w}$ is a $(1 \times p)$ vector times a $(p \times p)$ matrix times a $(p \times 1)$ vector. Thus the result is a (1×1) vector or a scalar which is what we were hoping for.

For the spread of the means or the distance between the means we can look at the quantity $\|\mathbb{E}(\mathbf{w}^\top \mathbf{X}) - \mathbb{E}(\mathbf{w}^\top \mathbf{Y})\|_2$ where \mathbf{X} are the random vectors in class 1 and \mathbf{Y} are the random vectors in class 2. The squared distance $\|\mathbb{E}(\mathbf{w}^\top \mathbf{X}) - \mathbb{E}(\mathbf{w}^\top \mathbf{Y})\|_2^2$ can be written as

$$\begin{aligned}
\underline{\|\mathbb{E}(\mathbf{w}^\top \mathbf{X}) - \mathbb{E}(\mathbf{w}^\top \mathbf{Y})\|_2^2} &= \langle \mathbb{E}(\mathbf{w}^\top \mathbf{X}) - \mathbb{E}(\mathbf{w}^\top \mathbf{Y}), \mathbb{E}(\mathbf{w}^\top \mathbf{X}) - \mathbb{E}(\mathbf{w}^\top \mathbf{Y}) \rangle \\
&= (\mathbb{E}[\mathbf{w}^\top \mathbf{X}] - \mathbb{E}[\mathbf{w}^\top \mathbf{Y}])^\top (\mathbb{E}[\mathbf{w}^\top \mathbf{X}] - \mathbb{E}[\mathbf{w}^\top \mathbf{Y}]) \\
&= (\mathbf{w}^\top \mathbb{E}[\mathbf{X}] - \mathbf{w}^\top \mathbb{E}[\mathbf{Y}])^\top (\mathbf{w}^\top \mathbb{E}[\mathbf{X}] - \mathbf{w}^\top \mathbb{E}[\mathbf{Y}]) \\
&= (\mathbb{E}[\mathbf{X}] - \mathbb{E}[\mathbf{Y}])^\top \mathbf{w} \mathbf{w}^\top (\mathbb{E}[\mathbf{X}] - \mathbb{E}[\mathbf{Y}])
\end{aligned}$$

which is a (1×1) vector so its transpose is the same as itself, so

$$= \mathbf{w}^\top (\mathbb{E}[\mathbf{X}] - \mathbb{E}[\mathbf{Y}]) (\mathbb{E}[\mathbf{X}] - \mathbb{E}[\mathbf{Y}])^\top \mathbf{w} = \underline{\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}}.$$

Since we want to maximize the distance between the projected means while minimizing the covariances of the projected variances this means we want to maximize $\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}$ while minimizing $\mathbf{w}^\top \Sigma_{\mathbf{X}\mathbf{X}} \mathbf{w} + \mathbf{w}^\top \Sigma_{\mathbf{Y}\mathbf{Y}} \mathbf{w} = \mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}$ with respect to \mathbf{w} . There are two natural ways to do this; either we can maximize the (generalized Rayleigh) quotient $\frac{\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}}{\mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}}$ or we can maximize the difference $\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w} - \mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}$. We will soon see that these two are equivalent in some way but first we need to add a condition to the optimization problem so that it is well defined: Consider the first optimization problem. For a given \mathbf{w} notice that stretching it by the scalar a won't have an impact on the value of the objective function $\frac{\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}}{\mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}}$ since $\frac{a \mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} a \mathbf{w}}{a \mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) a \mathbf{w}} = \frac{a^2 \mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}}{a^2 \mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}} = \frac{\mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w}}{\mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w}}$. Thus only the direction of the vector \mathbf{w} matters and we can add the condition $\mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w} = 1$, fixing the denominator at the same time.

The optimization problem that is ultimately solved is then the following:

$$\begin{aligned}
&\max_{\mathbf{w}} \quad \mathbf{w}^\top S_{\mathbf{X}\mathbf{Y}} \mathbf{w} \\
&\text{such that} \quad \mathbf{w}^\top (\Sigma_{\mathbf{X}\mathbf{X}} + \Sigma_{\mathbf{Y}\mathbf{Y}}) \mathbf{w} = 1.
\end{aligned} \tag{1}$$

This is readily solved using the Lagrange multiplier method. The Lagrangian is given by $\mathcal{L}(\mathbf{w}, \lambda) = \mathbf{w}^\top S_{\mathbf{XY}} \mathbf{w} - \lambda (\mathbf{w}^\top (\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}}) \mathbf{w} - 1)$ with the derivative

$$\frac{d\mathcal{L}(\mathbf{w}, \lambda)}{d\mathbf{w}} = 2S_{\mathbf{XY}} \mathbf{w} - 2\lambda (\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}}) \mathbf{w}.$$

For $\lambda = 1$ this is the same derivative as the derivative of the second optimization problem considered (with the difference between the projected spread and variances). Setting the derivative equal to zero gives the following relation:

$$S_{\mathbf{XY}} \mathbf{w} = \lambda (\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}}) \mathbf{w}.$$

This is called a generalized eigenvalue problem which can be changed into a regular eigenvector problem if $(\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}})$ is non-singular. If this is the case then we can make further progress.

Since $S_{\mathbf{XY}}$ is a product of two vectors it has rank 1 (the basis is one of the vectors and the weights are given by the other vector). Thus the result of the matrix multiplication $(\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}})^{-1} S_{\mathbf{XY}}$ also has at most rank 1 and thus has at most one eigenvalue not equal to one. The solution to the optimization problem is then the eigenvector of the matrix $(\Sigma_{\mathbf{XX}} + \Sigma_{\mathbf{YY}})^{-1} S_{\mathbf{XY}}$.

The algorithm usually referred to as linear discriminant analysis can be seen as first finding this projection and then choosing a point on the one-dimensional projection to use as the threshold for classification. This can be done by assuming things about the distribution of the data, for example assuming that the two classes are normally distributed around separate means but with the same covariance matrix. Using the assumptions and incorporating weighting based on the number of observations in each class one can determine a threshold such that the resulting classification is equivalent to the “normal” linear discriminant analysis algorithm.

3 A note on the solution

Usually, the covariance and distance matrices $\Sigma_{\mathbf{XX}}$, $\Sigma_{\mathbf{YY}}$ and $S_{\mathbf{XY}}$ are unknown and thus need to be estimated from some observed data. In other words the matrices depend on the observed data and one can even see that they lie in the space spanned by the observations in some sense.

Consider the maximum likelihood estimate of the expected value of class one ($\mathbb{E}(\mathbf{X})$) which is equal to

$$\mathbf{m}_{\mathbf{X}} := \frac{1}{N_{\mathbf{X}}} \sum_{i=1}^{N_{\mathbf{X}}} \mathbf{X}_i$$

which is a linear combination of the $N_{\mathbf{X}}$ observations \mathbf{X}_i and thus lies in the span of the observations. The analogous is true for $\mathbf{m}_{\mathbf{Y}}$, the maximum likelihood estimate of the expected value of class 2.

The covariances can be estimated as

$$\begin{aligned}\hat{\Sigma}_{\mathbf{X}\mathbf{X}} &:= \frac{1}{N_{\mathbf{X}}} \sum_{i=1}^{N_{\mathbf{X}}} (\mathbf{X}_i - \mathbf{m}_{\mathbf{X}}) (\mathbf{X}_i - \mathbf{m}_{\mathbf{X}})^{\top} \quad \text{and} \\ \hat{\Sigma}_{\mathbf{Y}\mathbf{Y}} &:= \frac{1}{N_{\mathbf{Y}}} \sum_{i=1}^{N_{\mathbf{Y}}} (\mathbf{Y}_i - \mathbf{m}_{\mathbf{Y}}) (\mathbf{Y}_i - \mathbf{m}_{\mathbf{Y}})^{\top},\end{aligned}$$

and the distance between the means as

$$\hat{S}_{\mathbf{X}\mathbf{Y}} := (\mathbf{m}_{\mathbf{X}} - \mathbf{m}_{\mathbf{Y}}) (\mathbf{m}_{\mathbf{X}} - \mathbf{m}_{\mathbf{Y}})^{\top}.$$

All three of these are matrices where the rows/columns are linear combinations of the observations \mathbf{X}_i and \mathbf{Y}_i .

For the solution to the optimization problem, the eigenvector, the same has to be true, i.e. the eigenvector has to be some linear combination of the observations. In other words,

$$\mathbf{w} = \sum_{i=1}^{N_{\mathbf{X}}} \alpha_{\mathbf{X}_i} \mathbf{X}_i + \sum_{i=1}^{N_{\mathbf{Y}}} \alpha_{\mathbf{Y}_i} \mathbf{Y}_i$$

or, by defining \mathbf{Z} as the matrix containing all $N_{\mathbf{X}} + N_{\mathbf{Y}} = N$ observations,

$$= \sum_{i=1}^N \alpha_i \mathbf{Z}_i.$$

Consider now a (possibly non-linear) transformation $h : \mathbb{R}^p \mapsto \mathbb{R}^P$. Suppose we wish to apply the (Fisher) linear discriminant analysis algorithm to this transformed dataset. Defining the vectors

$$\begin{aligned}\mathbf{m}_{h(\mathbf{X})} &:= \frac{1}{N_{\mathbf{X}}} \sum_{i=1}^{N_{\mathbf{X}}} h(\mathbf{X}_i) \quad \text{and} \\ \mathbf{m}_{h(\mathbf{Y})} &:= \frac{1}{N_{\mathbf{Y}}} \sum_{i=1}^{N_{\mathbf{Y}}} h(\mathbf{Y}_i),\end{aligned}$$

as well as the matrices

$$\begin{aligned}\hat{\Sigma}_{h(\mathbf{X})h(\mathbf{X})} &:= \frac{1}{N_{\mathbf{X}}} \sum_{i=1}^{N_{\mathbf{X}}} \left(h(\mathbf{X}_i) - \mathbf{m}_{h(\mathbf{X})} \right) \left(h(\mathbf{X}_i) - \mathbf{m}_{h(\mathbf{X})} \right)^{\top}, \\ \hat{\Sigma}_{h(\mathbf{Y})h(\mathbf{Y})} &:= \frac{1}{N_{\mathbf{Y}}} \sum_{i=1}^{N_{\mathbf{Y}}} \left(h(\mathbf{Y}_i) - \mathbf{m}_{h(\mathbf{Y})} \right) \left(h(\mathbf{Y}_i) - \mathbf{m}_{h(\mathbf{Y})} \right)^{\top} \quad \text{and} \\ \hat{S}_{h(\mathbf{X})h(\mathbf{Y})} &:= \left(\mathbf{m}_{h(\mathbf{X})} - \mathbf{m}_{h(\mathbf{Y})} \right) \left(\mathbf{m}_{h(\mathbf{X})} - \mathbf{m}_{h(\mathbf{Y})} \right)^{\top}\end{aligned}$$

we can try solving the optimization problem

$$\begin{aligned} \max_{\mathbf{w}} \quad & \mathbf{w}^\top \hat{S}_{h(\mathbf{X})h(\mathbf{Y})} \mathbf{w} \\ \text{such that} \quad & \mathbf{w}^\top \left(\hat{S}_{h(\mathbf{X})h(\mathbf{X})} + \hat{S}_{h(\mathbf{Y})h(\mathbf{Y})} \right) \mathbf{w} = 1, \end{aligned} \quad (2)$$

where \mathbf{w} is now a $(1 \times P)$ vector. If P is large this might be difficult, especially if $P > N$ since the matrix $\left(\hat{S}_{h(\mathbf{X})h(\mathbf{X})} + \hat{S}_{h(\mathbf{Y})h(\mathbf{Y})} \right) =: W_h$ will be singular and a similar approach as the solution to (1) would not work.

However, what we do know is that if there is a solution \mathbf{w} , it has to be of the form

$$\mathbf{w} = \sum_{i=1}^N \alpha_i h(\mathbf{Z}_i).$$

The linear combination of the observed transformed average, $\mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})}$, then has to have the form

$$\begin{aligned} \mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})} &= \frac{1}{N_{\mathbf{X}}} \mathbf{w}^\top \sum_{i=1}^{N_{\mathbf{X}}} h(\mathbf{X}_i) = \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^N \alpha_j h(\mathbf{Z}_j)^\top \sum_{i=1}^{N_{\mathbf{X}}} h(\mathbf{X}_i) \\ &= \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^N \sum_{i=1}^{N_{\mathbf{X}}} \alpha_j h(\mathbf{Z}_j)^\top h(\mathbf{X}_i) = \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^N \sum_{i=1}^{N_{\mathbf{X}}} \alpha_j \langle h(\mathbf{Z}_j), h(\mathbf{X}_i) \rangle. \end{aligned}$$

At this point we recognize the inner product between two high-dimensional vectors, $\langle h(\mathbf{Z}_j), h(\mathbf{X}_i) \rangle$, and realize that this can be computed using kernels.

4 Reformulating Fishers Linear Discriminant using Kernels

Kernels are functions that efficiently compute inner products in some high-dimensional (even infinite-dimensional) space. An advantage with using kernels is that when we have some mapping h to a higher dimensional space, the kernel doesn't need to explicitly compute the mapping but instead it implicitly computes in a high-dimensional space (given that the mapping h is such that it gives rise to a kernel). One can also go the other way around and fix a kernel that then determines what mapping one is working in. In general, kernels can always be used instead of inner products.

Continuing with the manipulation of $\mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})}$, the inner product $\langle h(\mathbf{Z}_j), h(\mathbf{X}_i) \rangle$ can be replaced by a suitable kernel $k(\mathbf{Z}_j, \mathbf{X}_i)$, meaning that

$$\underline{\mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})}} = \frac{1}{N_{\mathbf{X}}} \sum_{j=1}^N \sum_{i=1}^{N_{\mathbf{X}}} \alpha_j k(\mathbf{Z}_j, \mathbf{X}_i)$$

which can be seen as the multiplication

$$= \underline{\boldsymbol{\alpha}^\top \mathbf{M}_{\mathbf{X}}},$$

where α is an $(1 \times N)$ vector of the scalars $\alpha_j, j = 1, \dots, N$ and \mathbf{M}_X is the $(1 \times N)$ vector of the scalars $\frac{1}{N_X} \sum_{i=1}^{N_X} k(\mathbf{Z}_j, \mathbf{X}_i), j = 1, \dots, N$. Similarly we get that $\mathbf{w}^\top \mathbf{m}_{h(\mathbf{Y})} = \alpha^\top \mathbf{M}_Y$.

The spread between the two transformed and projected means is then given by

$$\begin{aligned} \left\| \mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})} - \mathbf{w}^\top \mathbf{m}_{h(\mathbf{Y})} \right\|_2^2 &= \left\| \alpha^\top \mathbf{M}_X - \alpha^\top \mathbf{M}_Y \right\|_2^2 = \langle \alpha^\top \mathbf{M}_X - \alpha^\top \mathbf{M}_Y, \alpha^\top \mathbf{M}_X - \alpha^\top \mathbf{M}_Y \rangle \\ &= (\alpha^\top \mathbf{M}_X - \alpha^\top \mathbf{M}_Y)^\top (\alpha^\top \mathbf{M}_X - \alpha^\top \mathbf{M}_Y) \\ &= (\mathbf{M}_X - \mathbf{M}_Y)^\top \alpha \alpha^\top (\mathbf{M}_X - \mathbf{M}_Y) \\ &= \alpha^\top (\mathbf{M}_X - \mathbf{M}_Y) (\mathbf{M}_X - \mathbf{M}_Y)^\top \alpha = \underline{\alpha^\top S_k \alpha}, \end{aligned}$$

this is the same calculation as the one performed for the original space and S_k is the $(N \times N)$ matrix defined as $S_k := (\mathbf{M}_X - \mathbf{M}_Y) (\mathbf{M}_X - \mathbf{M}_Y)^\top$.

For the covariance/variance within the groups we first need to take a look at the projection of a transformed variable,

$$\begin{aligned} \underline{\mathbf{w}^\top h(\mathbf{X})} &= \sum_{j=1}^N \alpha_j h(\mathbf{Z}_j)^\top h(\mathbf{X}) = \sum_{j=1}^N \alpha_j \langle h(\mathbf{Z}_j), h(\mathbf{X}) \rangle \\ &= \sum_{j=1}^N \alpha_j k(\mathbf{Z}_j, \mathbf{X}) = \underline{\alpha^\top \mathbf{K}_X}. \end{aligned}$$

The covariance $\underline{\text{Cov}(\mathbf{w}^\top h(\mathbf{X}), \mathbf{w}^\top h(\mathbf{X}))}$ is then estimated by

$$\begin{aligned} &\frac{1}{N_x} \sum_{i=1}^{N_x} \left((\mathbf{w}^\top h(\mathbf{X}_i) - \mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})}) (\mathbf{w}^\top h(\mathbf{X}_i) - \mathbf{w}^\top \mathbf{m}_{h(\mathbf{X})})^\top \right) \\ &= \frac{1}{N_x} \sum_{i=1}^{N_x} ((\alpha^\top \mathbf{K}_{X_i} - \alpha^\top \mathbf{M}_X) (\alpha^\top \mathbf{K}_{X_i} - \alpha^\top \mathbf{M}_X)^\top) \\ &= \frac{1}{N_x} \sum_{i=1}^{N_x} (\alpha^\top (\mathbf{K}_{X_i} - \mathbf{M}_X) (\mathbf{K}_{X_i} - \mathbf{M}_X)^\top \alpha) \\ &= \alpha^\top \left(\frac{1}{N_X} \sum_{i=1}^{N_X} (\mathbf{K}_{X_i} - \mathbf{M}_X) (\mathbf{K}_{X_i} - \mathbf{M}_X)^\top \right) \alpha = \underline{\alpha^\top \Sigma_{kX} \alpha} \end{aligned}$$

and similarly for class \mathbf{Y} we get the estimate $\alpha^\top \Sigma_{kY} \alpha$, where the $(N \times N)$ matrix describes the average observed distances from the transformed mean of the observations and is given by $\Sigma_{kY} := \left(\frac{1}{N_Y} \sum_{i=1}^{N_Y} (\mathbf{K}_{Y_i} - \mathbf{M}_Y) (\mathbf{K}_{Y_i} - \mathbf{M}_Y)^\top \right)$ (there might be more efficient ways to compute this). The two covariances can be summed, we define this as the matrix $W_k := \Sigma_{kX} + \Sigma_{kY}$.

In this way we have transformed the optimization problem (2) into the following optimization problem:

$$\begin{aligned} &\max_{\alpha} \quad \alpha^\top S_k \alpha \\ &\text{such that} \quad \alpha^\top W_k \alpha = 1, \end{aligned} \tag{3}$$

a problem where we already know that the solution α is given by the only eigenvector of the matrix $W_k^{-1} S_k$.

The projection of a new observation \mathbf{X} is given by $\mathbf{w}^\top h(\mathbf{X}) = \boldsymbol{\alpha} \mathbf{K}_{\mathbf{X}}$. The kernel k or the non-linear transformation h should be chosen using cross-validation. Classification can be done in the projected space using almost any classification method.

4.1 Some notes

In general, the matrix W_k might not be invertible according to [1] since N dimensional covariances are estimated using N observations. The authors suggest replacing W_k by $W_\gamma = W_k + \gamma I$ where I is the $(N \times N)$ identity matrix. For a sufficiently large γ , the resulting matrix will be invertible. This also has the benefit of having a regularizing effect, which is good since the high-dimensional space that the kernel allows for might lead to overfitting.

For the binary classification problem one can solve for the direction of the eigenvector by noting that $\boldsymbol{\alpha} \propto W_k^{-1} (\mathbf{M}_{\mathbf{X}} - \mathbf{M}_{\mathbf{Y}})$. Once the direction is known the magnitude does not really matter.

The method can be extended to multiple-classes and then gives as many eigenvectors or directions/projections as the number of classes minus 1. Then the estimation of the covariance matrices can be troublesome since all the combinations of the classes have to be taken into account.

5 Conclusion

Three optimization problems were presented and two of them solved in order to explain how the dimensionality reduction step of the linear discriminant classifier can be generalized to a kernelized version. The main steps were to formulate an optimization problem so that the resulting solution is pleasing, then to solve the optimization problem using Lagrange multipliers and to recognize the generalized eigenvalue problem. After the original problem is solved some knowledge of linear algebra and reproducing kernels is used to rework the problem into one where the data is only present in inner products so that they can be switched for kernels. After this one can implicitly work in infinite-dimensional spaces to get a good projection that can be used for classification.

References

- [1] S. Mika, G. Rätsch, J. Weston, B. Schölkopf and K.-R. Müller, *Fisher Discriminant Analysis With Kernels*. Neural Networks for Signal Processing IX: Proceedings of the 1999 IEEE Signal Processing Society Workshop, 1999.
- [2] T. Hastie, R. Tibshirani and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer, Second Edition, 2009.