

Fisher and Kernel Fisher Discriminant Analysis: Tutorial

Benyamin Ghogh

BGHOJOGH@UWATERLOO.CA

Department of Electrical and Computer Engineering,
Machine Learning Laboratory, University of Waterloo, Waterloo, ON, Canada

Fakhri Karray

KARRAY@UWATERLOO.CA

Department of Electrical and Computer Engineering,
Centre for Pattern Analysis and Machine Intelligence, University of Waterloo, Waterloo, ON, Canada

Mark Crowley

MCROWLEY@UWATERLOO.CA

Department of Electrical and Computer Engineering,
Machine Learning Laboratory, University of Waterloo, Waterloo, ON, Canada

Abstract

This is a detailed tutorial paper which explains the Fisher discriminant Analysis (FDA) and kernel FDA. We start with projection and reconstruction. Then, one- and multi-dimensional FDA subspaces are covered. Scatters in two- and then multi-classes are explained in FDA. Then, we discuss on the rank of the scatters and the dimensionality of the subspace. A real-life example is also provided for interpreting FDA. Then, possible singularity of the scatter is discussed to introduce robust FDA. PCA and FDA directions are also compared. We also prove that FDA and linear discriminant analysis are equivalent. Fisher forest is also introduced as an ensemble of fisher subspaces useful for handling data with different features and dimensionality. Afterwards, kernel FDA is explained for both one- and multi-dimensional subspaces with both two- and multi-classes. Finally, some simulations are performed on AT&T face dataset to illustrate FDA and compare it with PCA.

1. Introduction

Assume we have a dataset of *instances* or *data points* $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ with sample size n and dimensionality $\mathbf{x}_i \in \mathbb{R}^d$ and $\mathbf{y}_i \in \mathbb{R}^\ell$. The $\{\mathbf{x}_i\}_{i=1}^n$ are the input data to the model and the $\{\mathbf{y}_i\}_{i=1}^n$ are the observations (labels). We define $\mathbb{R}^{d \times n} \ni \mathbf{X} := [\mathbf{x}_1, \dots, \mathbf{x}_n]$ and $\mathbb{R}^{\ell \times n} \ni \mathbf{Y} := [\mathbf{y}_1, \dots, \mathbf{y}_n]$. We can also have an out-of-sample

data point, $\mathbf{x}_t \in \mathbb{R}^d$, which is not in the training set. If there are n_t out-of-sample data points, $\{\mathbf{x}_{t,i}\}_{i=1}^{n_t}$, we define $\mathbb{R}^{d \times n_t} \ni \mathbf{X}_t := [\mathbf{x}_{t,1}, \dots, \mathbf{x}_{t,n_t}]$. Usually, the data points exist on a subspace or sub-manifold. Subspace or manifold learning tries to learn this sub-manifold (Ghogh et al., 2019b).

Here, we consider the case where the observations $\{\mathbf{y}_i\}_{i=1}^n$ come from a discrete set so that the task is *classification*. Assume the dataset consists of c classes, $\{\mathbf{x}_i^{(1)}\}_{i=1}^{n_1}, \dots, \{\mathbf{x}_i^{(c)}\}_{i=1}^{n_c}$ where n_j denotes the sample size (cardinality) of the j -th class.

We want to find a subspace (or sub-manifold) which separates the classes as much as possible while the data also become as spread as possible. Fisher Discriminant Analysis (FDA) (Friedman et al., 2009) pursues this goal. It was first proposed in (Fisher, 1936) by Sir. Ronald Aylmer Fisher (1890 – 1962) who was a genius in statistics. He proposed many important concepts in the modern statistics, such as variance (Fisher, 1919), FDA (Fisher, 1936), Fisher information (Frieden, 2004), Analysis of Variance (ANOVA) (Fisher, 1992), etc. The paper (Fisher, 1936), which proposed FDA, was the first paper introducing the well-known Iris flower dataset. Note that Fisher’s work was mostly concentrating on the statistics in the area of genetics. Much of his work was about variance making no wonder for us why FDA is all about variance and scatters.

Kernel FDA (Mika et al., 1999; 2000) performs the goal of FDA in the feature space. The FDA and kernel FDA have had many different applications. Some examples for applications of FDA are face recognition (Fisherfaces) (Belhumeur et al., 1997; Etemad & Chellappa, 1997; Zhao et al., 1999), action recognition (Fisherposes) (Ghogh et al., 2017; Mokari et al., 2018), and gesture recognition (Samadani et al., 2013). Some examples for applications

of kernel FDA are face recognition (kernel Fisherfaces) (Yang, 2002; Liu et al., 2004) and palmprint Recognition (Wang & Ruan, 2006).

In the literature, sometimes, FDA is referred to as Linear Discriminant Analysis (LDA) or Fisher LDA (FLDA). This is because FDA and LDA (Ghojogh & Crowley, 2019a) are equivalent although LDA is a classification method and not a subspace learning algorithm. In this paper, we will prove why they are equivalent.

2. Projection Formulation

2.1. Projection

Assume we have a data point $\mathbf{x} \in \mathbb{R}^d$. We want to project this data point onto the vector space spanned by p vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ where each vector is d -dimensional and usually $p \ll d$. We stack these vectors column-wise in matrix $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p] \in \mathbb{R}^{d \times p}$. In other words, we want to project \mathbf{x} onto the column space of \mathbf{U} , denoted by $\text{Col}(\mathbf{U})$. The projection of $\mathbf{x} \in \mathbb{R}^d$ onto $\text{Col}(\mathbf{U}) \in \mathbb{R}^p$ and then its representation in the \mathbb{R}^d (its reconstruction) can be seen as a linear system of equations:

$$\mathbb{R}^d \ni \hat{\mathbf{x}} := \mathbf{U}\boldsymbol{\beta}, \quad (1)$$

where we should find the unknown coefficients $\boldsymbol{\beta} \in \mathbb{R}^p$.

If the \mathbf{x} lies in the $\text{Col}(\mathbf{U})$ or $\text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$, this linear system has exact solution, so $\hat{\mathbf{x}} = \mathbf{x} = \mathbf{U}\boldsymbol{\beta}$. However, if \mathbf{x} does not lie in this space, there is no any solution $\boldsymbol{\beta}$ for $\mathbf{x} = \mathbf{U}\boldsymbol{\beta}$ and we should solve for projection of \mathbf{x} onto $\text{Col}(\mathbf{U})$ or $\text{span}\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ and then its reconstruction. In other words, we should solve for Eq. (1). In this case, $\hat{\mathbf{x}}$ and \mathbf{x} are different and we have a residual:

$$\mathbf{r} = \mathbf{x} - \hat{\mathbf{x}} = \mathbf{x} - \mathbf{U}\boldsymbol{\beta}, \quad (2)$$

which we want to be small. As can be seen in Fig. 1, the smallest residual vector is orthogonal to $\text{Col}(\mathbf{U})$; therefore:

$$\begin{aligned} \mathbf{x} - \mathbf{U}\boldsymbol{\beta} \perp \mathbf{U} &\implies \mathbf{U}^\top (\mathbf{x} - \mathbf{U}\boldsymbol{\beta}) = 0, \\ &\implies \boldsymbol{\beta} = (\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{x}. \end{aligned} \quad (3)$$

It is noteworthy that the Eq. (3) is also the formula of coefficients in linear regression (Friedman et al., 2009) where the input data are the rows of \mathbf{U} and the labels are \mathbf{x} ; however, our goal here is different.

Plugging Eq. (3) in Eq. (1) gives us:

$$\hat{\mathbf{x}} = \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top \mathbf{x}.$$

We define:

$$\mathbb{R}^{d \times d} \ni \boldsymbol{\Pi} := \mathbf{U}(\mathbf{U}^\top \mathbf{U})^{-1} \mathbf{U}^\top, \quad (4)$$

as ‘‘projection matrix’’ because it projects \mathbf{x} onto $\text{Col}(\mathbf{U})$ (and reconstructs back). Note that $\boldsymbol{\Pi}$ is also referred to as

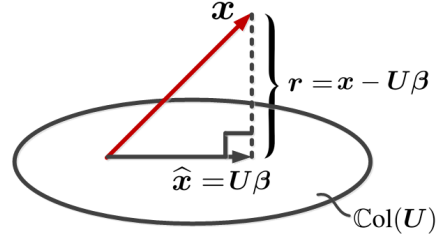


Figure 1. The residual and projection onto the column space of \mathbf{U} .

the ‘‘hat matrix’’ in the literature because it puts a hat on top of \mathbf{x} .

If the vectors $\{\mathbf{u}_1, \dots, \mathbf{u}_p\}$ are orthonormal (the matrix \mathbf{U} is orthogonal), we have $\mathbf{U}^\top = \mathbf{U}^{-1}$ and thus $\mathbf{U}^\top \mathbf{U} = \mathbf{I}$. Therefore, Eq. (4) is simplified:

$$\boldsymbol{\Pi} = \mathbf{U}\mathbf{U}^\top. \quad (5)$$

So, we have:

$$\hat{\mathbf{x}} = \boldsymbol{\Pi} \mathbf{x} = \mathbf{U}\mathbf{U}^\top \mathbf{x}. \quad (6)$$

2.2. Projection onto a Subspace

In subspace learning, the projection of a vector $\mathbf{x} \in \mathbb{R}^d$ onto the column space of $\mathbf{U} \in \mathbb{R}^{d \times p}$ (a p -dimensional subspace spanned by $\{\mathbf{u}_j\}_{j=1}^p$ where $\mathbf{u}_j \in \mathbb{R}^d$) is defined as:

$$\mathbb{R}^p \ni \tilde{\mathbf{x}} := \mathbf{U}^\top \mathbf{x}, \quad (7)$$

$$\mathbb{R}^d \ni \hat{\mathbf{x}} := \mathbf{U}\mathbf{U}^\top \mathbf{x} = \mathbf{U}\tilde{\mathbf{x}}, \quad (8)$$

where $\tilde{\mathbf{x}}$ and $\hat{\mathbf{x}}$ denote the projection and reconstruction of \mathbf{x} , respectively.

If we have n data points, $\{\mathbf{x}_i\}_{i=1}^n$, which can be stored column-wise in a matrix $\mathbf{X} \in \mathbb{R}^{d \times n}$, the projection and reconstruction of \mathbf{X} are defined as:

$$\mathbb{R}^{p \times n} \ni \tilde{\mathbf{X}} := \mathbf{U}^\top \mathbf{X}, \quad (9)$$

$$\mathbb{R}^{d \times n} \ni \hat{\mathbf{X}} := \mathbf{U}\mathbf{U}^\top \mathbf{X} = \mathbf{U}\tilde{\mathbf{X}}, \quad (10)$$

respectively.

If we have an out-of-sample data point \mathbf{x}_t which was not used in calculation of \mathbf{U} , the projection and reconstruction of it are defined as:

$$\mathbb{R}^p \ni \tilde{\mathbf{x}}_t := \mathbf{U}^\top \mathbf{x}_t, \quad (11)$$

$$\mathbb{R}^d \ni \hat{\mathbf{x}}_t := \mathbf{U}\mathbf{U}^\top \mathbf{x}_t = \mathbf{U}\tilde{\mathbf{x}}_t, \quad (12)$$

respectively.

In case we have n_t out-of-sample data points, $\{\mathbf{x}_{t,i}\}_{i=1}^{n_t}$, which can be stored column-wise in a matrix $\mathbf{X}_t \in \mathbb{R}^{d \times n_t}$, the projection and reconstruction of \mathbf{X}_t are defined as:

$$\mathbb{R}^{p \times n_t} \ni \tilde{\mathbf{X}}_t := \mathbf{U}^\top \mathbf{X}_t, \quad (13)$$

$$\mathbb{R}^{d \times n_t} \ni \hat{\mathbf{X}}_t := \mathbf{U}\mathbf{U}^\top \mathbf{X}_t = \mathbf{U}\tilde{\mathbf{X}}_t, \quad (14)$$

respectively.

For the properties of the projection matrix U , refer to (Ghojogh & Crowley, 2019c).

2.2.1. PROJECTION ONTO A ONE-DIMENSIONAL SUBSPACE

Considering the data $\{x_i\}_{i=1}^n$, the mean of data is:

$$\mathbb{R}^d \ni \mu_x := \frac{1}{n} \sum_{i=1}^n x_i, \quad (15)$$

and the centered data point x is:

$$\mathbb{R}^d \ni \check{x} := x - \mu_x. \quad (16)$$

The centered data X is:

$$\mathbb{R}^{d \times n} \ni \check{X} := X - \mu_x = XH, \quad (17)$$

where $\check{X} = [\check{x}_1, \dots, \check{x}_n] \in \mathbb{R}^{d \times n}$ and $\mathbb{R}^{n \times n} \ni H := I - (1/n)\mathbf{1}\mathbf{1}^\top$ is the centering matrix (see Appendix A in (Ghojogh & Crowley, 2019c)).

In Eq. (8), if $p = 1$, we are projecting x onto only one vector u and reconstruct it. If the data point is centered, the reconstruction is:

$$\hat{x} = uu^\top \check{x}.$$

The squared length (squared ℓ_2 -norm) of this reconstructed vector is:

$$\begin{aligned} \|\hat{x}\|_2^2 &= \|uu^\top \check{x}\|_2^2 = (uu^\top \check{x})^\top (uu^\top \check{x}) \\ &= \check{x}^\top \underbrace{u u^\top u}_{1} u^\top \check{x} \stackrel{(a)}{=} \check{x}^\top u u^\top \check{x} \stackrel{(b)}{=} u^\top \check{x} \check{x}^\top u, \end{aligned} \quad (18)$$

where (a) is because u is a unit (normal) vector, i.e., $u^\top u = \|u\|_2^2 = 1$, and (b) is because $\check{x}^\top u = u^\top \check{x} \in \mathbb{R}$.

Suppose we have n data points $\{x_i\}_{i=1}^n$ where $\{\check{x}_i\}_{i=1}^n$ are the centered data. The summation of the squared lengths of their projections $\{\hat{x}_i\}_{i=1}^n$ is:

$$\sum_{i=1}^n \|\hat{x}_i\|_2^2 \stackrel{(18)}{=} \sum_{i=1}^n u^\top \check{x}_i \check{x}_i^\top u = u^\top \left(\sum_{i=1}^n \check{x}_i \check{x}_i^\top \right) u. \quad (19)$$

Considering $\check{X} = [\check{x}_1, \dots, \check{x}_n] \in \mathbb{R}^{d \times n}$, we have:

$$\begin{aligned} \mathbb{R}^{d \times d} \ni S &:= \sum_{i=1}^n (x_i - \mu_x)(x_i - \mu_x)^\top \stackrel{(16)}{=} \sum_{i=1}^n \check{x}_i \check{x}_i^\top \\ &= \check{X} \check{X}^\top \stackrel{(17)}{=} X H H X^\top, \end{aligned} \quad (20)$$

where S is called the ‘‘covariance matrix’’ or ‘‘scatter matrix’’. If the data were already centered, we would have $S = X X^\top$.

Plugging Eq. (20) in Eq. (19) gives us:

$$\sum_{i=1}^n \|\hat{x}_i\|_2^2 = u^\top S u. \quad (21)$$

Note that we can also say that $u^\top S u$ is the variance of the projected data onto PCA subspace. In other words, $u^\top S u = \text{Var}(u^\top \check{X})$. This makes sense because when some non-random thing (here u) is multiplied to the random data (here \check{X}), it will have squared (quadratic) effect on variance, and $u^\top S u$ is quadratic in u .

Therefore, $u^\top S u$ can be interpreted in two ways: (I) the squared length of reconstruction and (II) the variance of projection.

If we consider the n data points in the matrix $X \in \mathbb{R}^{d \times n}$, the squared length of reconstruction of the centered data is:

$$\begin{aligned} \|\widehat{X}\|_F^2 &= \|u u^\top \check{X}\|_F^2 = \text{tr}((u u^\top \check{X})^\top (u u^\top \check{X})) \\ &= \text{tr}(\check{X}^\top \underbrace{u u^\top u}_{1} \check{X}) \stackrel{(a)}{=} \text{tr}(\check{X}^\top u u^\top \check{X}) \\ &\stackrel{(b)}{=} \text{tr}(u^\top \check{X} \check{X}^\top u) \stackrel{(c)}{=} u^\top \check{X} \check{X}^\top u \stackrel{(20)}{=} u^\top S u, \end{aligned}$$

where $\text{tr}(\cdot)$ denotes the trace of matrix, (a) is because u is a unit vector, (b) is because of the cyclic property of the trace, and (c) is because $u^\top \check{X} \check{X}^\top u$ is a scalar. Hence, we have:

$$\|\widehat{X}\|_F^2 = u^\top S u. \quad (22)$$

2.2.2. PROJECTION ONTO A MULTI-DIMENSIONAL SUBSPACE

In Eq. (10), if $p > 1$, we are projecting the data onto a subspace with dimensionality more than one (spanned by $\{u_j\}_{j=1}^p$) and then reconstruct back. If the data X are assumed to be centered, the reconstruction is:

$$\widehat{X} = U U^\top \check{X}.$$

The squared length (squared Frobenius Norm) of this reconstructed matrix is:

$$\begin{aligned} \|\widehat{X}\|_F^2 &= \|U U^\top \check{X}\|_F^2 = \text{tr}((U U^\top \check{X})^\top (U U^\top \check{X})) \\ &= \text{tr}(\check{X}^\top \underbrace{U U^\top U}_{I} U^\top \check{X}) \stackrel{(a)}{=} \text{tr}(\check{X}^\top U U^\top \check{X}) \\ &\stackrel{(b)}{=} \text{tr}(U^\top \check{X} \check{X}^\top U) \stackrel{(20)}{=} \text{tr}(U^\top S U), \end{aligned}$$

where (a) is because U is an orthogonal matrix (its columns are orthonormal) and (b) is because of the cyclic property of trace. Thus, we have:

$$\|\widehat{X}\|_F^2 = \text{tr}(U^\top S U). \quad (23)$$

3. Fisher Discriminant Analysis

3.1. One-dimensional Subspace

3.1.1. SCATTERS IN TWO-CLASS CASE

Assume we have two classes, $\{\mathbf{x}_i^{(1)}\}_{i=1}^{n_1}$ and $\{\mathbf{x}_i^{(2)}\}_{i=1}^{n_2}$, where n_1 and n_2 denote the sample size of the first and second class, respectively, and $\mathbf{x}_i^{(j)}$ denotes the i -th instance of the j -th class.

If the data instances of the j -th class are projected onto a one-dimensional subspace (vector \mathbf{u}) by $\mathbf{u}^\top \mathbf{x}_i^{(j)}$, the mean and the variance of the projected data are $\mathbf{u}^\top \boldsymbol{\mu}_j$ and $\mathbf{u}^\top \mathbf{S}_j \mathbf{u}$, respectively, where $\boldsymbol{\mu}_j$ and \mathbf{S}_j are the mean and covariance matrix (scatter) of the j -th class. The mean of the j -th class is:

$$\mathbb{R}^d \ni \boldsymbol{\mu}_j := \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbf{x}_i^{(j)}. \quad (24)$$

According to Appendix A, after projection onto the one-dimensional subspace, the distance between the means of classes is:

$$\begin{aligned} \mathbb{R} \ni d_B &:= (\mathbf{u}^\top \boldsymbol{\mu}_1 - \mathbf{u}^\top \boldsymbol{\mu}_2)^\top (\mathbf{u}^\top \boldsymbol{\mu}_1 - \mathbf{u}^\top \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u} \mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &\stackrel{(a)}{=} \text{tr}((\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u} \mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)) \\ &\stackrel{(b)}{=} \text{tr}(\mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u}) \\ &\stackrel{(c)}{=} \mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u} \stackrel{(d)}{=} \mathbf{u}^\top \mathbf{S}_B \mathbf{u}, \end{aligned} \quad (25)$$

where (a) is because $(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u} \mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ is a scalar, (b) is because of the cyclic property of trace, (c) is because $\mathbf{u}^\top (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \mathbf{u}$ is a scalar, and (d) is because we define:

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_B := (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top, \quad (26)$$

as the *between-scatter* of classes.

The Eq. (25) can also be interpreted according to Eq. (22): the d_B is the variance of projection of the class means or the squared length of reconstruction of the class means.

We saw that the variance of projection is $\mathbf{u}^\top \mathbf{S}_j \mathbf{u}$ for the j -th class. If we add up the variances of projections of the two classes, we have:

$$\begin{aligned} \mathbb{R} \ni d_W &:= \mathbf{u}^\top \mathbf{S}_1 \mathbf{u} + \mathbf{u}^\top \mathbf{S}_2 \mathbf{u} = \mathbf{u}^\top (\mathbf{S}_1 + \mathbf{S}_2) \mathbf{u} \\ &\stackrel{(a)}{=} \mathbf{u}^\top \mathbf{S}_W \mathbf{u}, \end{aligned} \quad (27)$$

where:

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_W := \mathbf{S}_1 + \mathbf{S}_2, \quad (28)$$

is the *within-scatter* of classes. According to Eq. (22), the d_W is the summation of projection variance of class instances or the summation of the reconstruction length of class instances.

3.1.2. SCATTERS IN MULTI-CLASS CASE: VARIANT 1

Assume $\{\mathbf{x}_i^{(j)}\}_{i=1}^{n_j}$ are the instances of the j -th class where we have multiple number of classes. In this case, the *between-scatter* is defined as:

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_B := \sum_{j=1}^c (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^\top, \quad (29)$$

where c is the number of classes and:

$$\mathbb{R}^d \ni \boldsymbol{\mu} := \frac{1}{\sum_{k=1}^c n_k} \sum_{j=1}^c n_j \boldsymbol{\mu}_j = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i, \quad (30)$$

is the weighted mean of means of classes or the total mean of data.

It is noteworthy that some researches define the between-scatter in a weighted way:

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_B := \sum_{j=1}^c n_j (\boldsymbol{\mu}_j - \boldsymbol{\mu})(\boldsymbol{\mu}_j - \boldsymbol{\mu})^\top. \quad (31)$$

If we extend the Eq. (28) to c number of classes, the *within-scatter* is defined as:

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_W := \sum_{j=1}^c \mathbf{S}_j \quad (32)$$

$$\stackrel{(20)}{=} \sum_{j=1}^c \sum_{i=1}^{n_j} (\mathbf{x}_i^{(j)} - \boldsymbol{\mu}_j)(\mathbf{x}_i^{(j)} - \boldsymbol{\mu}_j)^\top, \quad (33)$$

where n_j is the sample size of the j -th class.

In this case, the d_B and d_W are:

$$\mathbb{R} \ni d_B := \mathbf{u}^\top \mathbf{S}_B \mathbf{u}, \quad (34)$$

$$\mathbb{R} \ni d_W := \mathbf{u}^\top \mathbf{S}_W \mathbf{u}, \quad (35)$$

where \mathbf{S}_B and \mathbf{S}_W are Eqs. (29) and (33).

3.1.3. SCATTERS IN MULTI-CLASS CASE: VARIANT 2

There is another variant for multi-class case in FDA. In this variant, the within-scatter is the same as Eq. (33). The between-scatter is, however, different.

The *total-scatter* is defined as the covariance matrix of the whole data, regardless of classes (Welling, 2005):

$$\mathbb{R}^{d \times d} \ni \mathbf{S}_T := \frac{1}{n} \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top, \quad (36)$$

where the total mean $\boldsymbol{\mu}$ is the Eq. (30). We can also use the scaled total-scatter by dropping the $1/n$ factor. On the other hand, the total scatter is equal to the summation of the within- and between-scatters:

$$\mathbf{S}_T = \mathbf{S}_W + \mathbf{S}_B. \quad (37)$$

Therefore, the between-scatter, in this variant, is obtained as:

$$\mathbf{S}_B := \mathbf{S}_T - \mathbf{S}_W. \quad (38)$$

3.1.4. FISHER SUBSPACE: VARIANT 1

In FDA, we want to maximize the projection variance (scatter) of means of classes and minimize the projection variance (scatter) of class instances. In other words, we want to maximize d_B and minimize d_W . The reason is that after projection, we want the within scatter of every class to be small and the between scatter of classes to be large; therefore, the instances of every class get close to one another and the classes get far from each other. The two mentioned optimization problems are:

$$\underset{\mathbf{u}}{\text{maximize}} \quad d_B(\mathbf{u}), \quad (39)$$

$$\underset{\mathbf{u}}{\text{minimize}} \quad d_W(\mathbf{u}). \quad (40)$$

We can merge these two optimization problems as a regularized optimization problem:

$$\underset{\mathbf{u}}{\text{maximize}} \quad d_B(\mathbf{u}) - \alpha d_W(\mathbf{u}), \quad (41)$$

where $\alpha > 0$ is the regularization parameter. Another way of merging Eqs. (39) and (40) is:

$$\underset{\mathbf{u}}{\text{maximize}} \quad f(\mathbf{u}) := \frac{d_B(\mathbf{u})}{d_W(\mathbf{u})} = \frac{\mathbf{u}^\top \mathbf{S}_B \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}}, \quad (42)$$

where $f(\mathbf{u}) \in \mathbb{R}$ is referred to as the *Fisher criterion* (Xu & Lu, 2006). The Fisher criterion is a generalized Rayleigh-Ritz Quotient (see Appendix B):

$$f(\mathbf{u}) \stackrel{(158)}{=} R(\mathbf{S}_B, \mathbf{S}_W; \mathbf{u}). \quad (43)$$

According to Eq. (165) in Appendix B, the optimization in Eq. (42) is equivalent to:

$$\begin{aligned} &\underset{\mathbf{u}}{\text{maximize}} \quad \mathbf{u}^\top \mathbf{S}_B \mathbf{u} \\ &\text{subject to} \quad \mathbf{u}^\top \mathbf{S}_W \mathbf{u} = 1. \end{aligned} \quad (44)$$

The Lagrangian (Boyd & Vandenberghe, 2004) is:

$$\mathcal{L} = \mathbf{w}^\top \mathbf{S}_B \mathbf{w} - \lambda(\mathbf{w}^\top \mathbf{S}_W \mathbf{w} - 1),$$

where λ is the Lagrange multiplier. Equating the derivative of \mathcal{L} to zero gives:

$$\begin{aligned} \mathbb{R}^d \ni \frac{\partial \mathcal{L}}{\partial \mathbf{u}} &= 2 \mathbf{S}_B \mathbf{u} - 2 \lambda \mathbf{S}_W \mathbf{u} \stackrel{\text{set}}{=} \mathbf{0} \\ \implies 2 \mathbf{S}_B \mathbf{u} &= 2 \lambda \mathbf{S}_W \mathbf{u} \implies \mathbf{S}_B \mathbf{u} = \lambda \mathbf{S}_W \mathbf{u}, \end{aligned} \quad (45)$$

which is a generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ according to (Ghojogh et al., 2019a). The \mathbf{u} is the eigenvector with the largest eigenvalue (because the optimization is maximization) and the λ is the corresponding eigenvalue. The \mathbf{u} is referred to as the *Fisher direction* or *Fisher axis*. The projection and reconstruction are according to Eqs. (9) and (10), respectively, where $\mathbf{u} \in \mathbb{R}^d$ is used instead of

$\mathbf{U} \in \mathbb{R}^{d \times p}$. The out-of-sample projection and reconstruction are according to Eqs. (13) and (14), respectively, with \mathbf{u} rather than \mathbf{U} .

One possible solution to the generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ is (Ghojogh et al., 2019a):

$$\begin{aligned} \mathbf{S}_B \mathbf{u} &= \lambda \mathbf{S}_W \mathbf{u} \implies \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{u} = \lambda \mathbf{u} \\ \implies \mathbf{u} &= \text{eig}(\mathbf{S}_W^{-1} \mathbf{S}_B), \end{aligned} \quad (46)$$

where $\text{eig}(\cdot)$ denotes the eigenvector of the matrix with the largest eigenvalue. Although the solution in Eq. (46) is a little dirty (Ghojogh et al., 2019a) because \mathbf{S}_W might be singular and not invertible, but this solution is very common for FDA. In some researches, the diagonal of \mathbf{S}_W is strengthened slightly to make it full rank and invertible (Ghojogh et al., 2019a):

$$\mathbf{u} = \text{eig}((\mathbf{S}_W + \varepsilon \mathbf{I})^{-1} \mathbf{S}_B), \quad (47)$$

where ε is a very small positive number, large enough to make \mathbf{S}_W full rank.

In a future section, we will cover robust FDA which tackles this problem. On the other hand, the generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ has a rigorous solution (Ghojogh et al., 2019a; Wang, 2015) which does not require non-singularity of \mathbf{S}_W .

Another way to solve the optimization in Eq. (42) is taking derivative from the Fisher criterion:

$$\begin{aligned} \mathbb{R}^d \ni \frac{\partial f(\mathbf{u})}{\partial \mathbf{u}} &= \frac{1}{(\mathbf{u}^\top \mathbf{S}_W \mathbf{u})^2} \times \\ &\left[(\mathbf{u}^\top \mathbf{S}_W \mathbf{u})(2 \mathbf{S}_B \mathbf{u}) - (\mathbf{u}^\top \mathbf{S}_B \mathbf{u})(2 \mathbf{S}_W \mathbf{u}) \right] \stackrel{\text{set}}{=} \mathbf{0} \\ \stackrel{(a)}{\implies} \mathbf{S}_B \mathbf{u} &= \frac{\mathbf{u}^\top \mathbf{S}_B \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}} \mathbf{S}_W \mathbf{u}, \end{aligned} \quad (48)$$

where (a) is because $\mathbf{u}^\top \mathbf{S}_W \mathbf{u}$ is a scalar. The Eq. (48) which is a generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ (Ghojogh et al., 2019a) with \mathbf{u} and $(\mathbf{u}^\top \mathbf{S}_B \mathbf{u})/(\mathbf{u}^\top \mathbf{S}_W \mathbf{u})$ as the eigenvector with the largest eigenvalue (because the optimization is maximization) and the corresponding eigenvalue, respectively. Therefore, the *Fisher criterion* is the eigenvalue of the Fisher direction.

3.1.5. FISHER SUBSPACE: VARIANT 2

Another way to find the FDA direction is to consider another version of Fisher criterion. According to Eq. (38) for \mathbf{S}_B , the Fisher criterion becomes (Welling, 2005):

$$\begin{aligned} f(\mathbf{u}) &= \frac{\mathbf{u}^\top \mathbf{S}_B \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}} \stackrel{(38)}{=} \frac{\mathbf{u}^\top (\mathbf{S}_T - \mathbf{S}_W) \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}} \\ &= \frac{\mathbf{u}^\top \mathbf{S}_T \mathbf{u} - \mathbf{u}^\top \mathbf{S}_W \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}} = \frac{\mathbf{u}^\top \mathbf{S}_T \mathbf{u}}{\mathbf{u}^\top \mathbf{S}_W \mathbf{u}} - 1. \end{aligned} \quad (49)$$

The -1 is a constant and is dropped in the optimization; therefore:

$$\begin{aligned} & \underset{\mathbf{u}}{\text{maximize}} \quad \mathbf{u}^\top \mathbf{S}_T \mathbf{u} \\ & \text{subject to} \quad \mathbf{u}^\top \mathbf{S}_W \mathbf{u} = 1, \end{aligned} \quad (50)$$

whose solution is similarly obtained as:

$$\mathbf{S}_T \mathbf{u} = \lambda \mathbf{S}_W \mathbf{u}, \quad (51)$$

which is a generalized eigenvalue problem $(\mathbf{S}_T, \mathbf{S}_W)$ according to (Ghojogh et al., 2019a).

3.2. Multi-dimensional Subspace

In case the Fisher subspace is the span of several Fisher directions, $\{\mathbf{u}_j\}_{j=1}^p$ where $\mathbf{u}_j \in \mathbb{R}^d$, the d_B and d_W are defined as:

$$\mathbb{R} \ni d_B := \text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U}), \quad (52)$$

$$\mathbb{R} \ni d_W := \text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U}), \quad (53)$$

where $\mathbb{R}^{d \times p} \ni \mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_p]$. In this case, maximizing the *Fisher criterion* is:

$$\underset{\mathbf{U}}{\text{maximize}} \quad f(\mathbf{U}) := \frac{d_B(\mathbf{U})}{d_W(\mathbf{U})} = \frac{\text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U})}{\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})}. \quad (54)$$

The Fisher criterion $f(\mathbf{U})$ is a generalized Rayleigh-Ritz Quotient (see Appendix B). According to Eq. (165) in Appendix B, the optimization in Eq. (54) is equivalent to:

$$\begin{aligned} & \underset{\mathbf{U}}{\text{maximize}} \quad \text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U}) \\ & \text{subject to} \quad \mathbf{U}^\top \mathbf{S}_W \mathbf{U} = \mathbf{I}. \end{aligned} \quad (55)$$

The Lagrangian (Boyd & Vandenberghe, 2004) is:

$$\mathcal{L} = \text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U}) - \text{tr}(\mathbf{\Lambda}^\top (\mathbf{U}^\top \mathbf{S}_W \mathbf{U} - \mathbf{I})),$$

where $\mathbf{\Lambda} \in \mathbb{R}^{d \times d}$ is a diagonal matrix whose diagonal entries are the Lagrange multipliers. Equating the derivative of \mathcal{L} to zero gives:

$$\begin{aligned} \mathbb{R}^{d \times p} \ni \frac{\partial \mathcal{L}}{\partial \mathbf{U}} &= 2 \mathbf{S}_B \mathbf{U} - 2 \mathbf{S}_W \mathbf{U} \mathbf{\Lambda} \stackrel{\text{set}}{=} \mathbf{0} \\ \implies 2 \mathbf{S}_B \mathbf{U} &= 2 \mathbf{S}_W \mathbf{U} \mathbf{\Lambda} \implies \mathbf{S}_B \mathbf{U} = \mathbf{S}_W \mathbf{U} \mathbf{\Lambda}, \end{aligned} \quad (56)$$

which is a generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ according to (Ghojogh et al., 2019a). The columns of \mathbf{U} are the eigenvectors sorted by largest to smallest eigenvalues (because the optimization is maximization) and the diagonal entries of $\mathbf{\Lambda}$ are the corresponding eigenvalues. The columns of \mathbf{U} are referred to as the *Fisher directions* or *Fisher axes*. The projection and reconstruction are according to Eqs. (9) and (10), respectively. The out-of-sample projection and reconstruction are according to Eqs. (13) and (14), respectively.

One possible solution to the generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ is (Ghojogh et al., 2019a):

$$\begin{aligned} \mathbf{S}_B \mathbf{U} &= \mathbf{S}_W \mathbf{U} \mathbf{\Lambda} \implies \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{U} = \mathbf{U} \mathbf{\Lambda} \\ \implies \mathbf{U} &= \text{eig}(\mathbf{S}_W^{-1} \mathbf{S}_B), \end{aligned} \quad (57)$$

where $\text{eig}(\cdot)$ denotes the eigenvectors of the matrix stacked column-wise. Again, we can have (Ghojogh et al., 2019a):

$$\mathbf{U} = \text{eig}((\mathbf{S}_W + \varepsilon \mathbf{I})^{-1} \mathbf{S}_B), \quad (58)$$

Another way to solve the optimization in Eq. (54) is taking derivative from the Fisher criterion:

$$\begin{aligned} \mathbb{R}^{d \times p} \ni \frac{\partial f(\mathbf{U})}{\partial \mathbf{U}} &= \frac{1}{(\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U}))^2} \times \\ & \left[\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})(2 \mathbf{S}_B \mathbf{U}) - \text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U})(2 \mathbf{S}_W \mathbf{U}) \right] \stackrel{\text{set}}{=} \mathbf{0} \\ \stackrel{(a)}{\implies} \mathbf{S}_B \mathbf{U} &= \frac{\text{tr}(\mathbf{U}^\top \mathbf{S}_B \mathbf{U})}{\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})} \mathbf{S}_W \mathbf{U}, \end{aligned} \quad (59)$$

where (a) is because $\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})$ is a scalar. The Eq. (59) which is a generalized eigenvalue problem $(\mathbf{S}_B, \mathbf{S}_W)$ (Ghojogh et al., 2019a) with columns of \mathbf{U} as the eigenvectors and $(\mathbf{u}_j^\top \mathbf{S}_B \mathbf{u}_j) / (\mathbf{u}_j^\top \mathbf{S}_W \mathbf{u}_j)$ as the j -th largest eigenvalue (because the optimization is maximization).

Again, another way to find the FDA directions is to consider another version of Fisher criterion. According to Eq. (38) for \mathbf{S}_B , the Fisher criterion becomes (Welling, 2005):

$$f(\mathbf{U}) = \frac{\text{tr}(\mathbf{U}^\top (\mathbf{S}_T - \mathbf{S}_W) \mathbf{U})}{\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})} = \frac{\text{tr}(\mathbf{U}^\top \mathbf{S}_T \mathbf{U})}{\text{tr}(\mathbf{U}^\top \mathbf{S}_W \mathbf{U})} - 1. \quad (60)$$

The -1 is a constant and is dropped in the optimization; therefore:

$$\begin{aligned} & \underset{\mathbf{U}}{\text{maximize}} \quad \text{tr}(\mathbf{U}^\top \mathbf{S}_T \mathbf{U}) \\ & \text{subject to} \quad \mathbf{U}^\top \mathbf{S}_W \mathbf{U} = \mathbf{I}, \end{aligned} \quad (61)$$

whose solution is similarly obtained as:

$$\mathbf{S}_T \mathbf{U} = \mathbf{S}_W \mathbf{U} \mathbf{\Lambda}, \quad (62)$$

which is a generalized eigenvalue problem $(\mathbf{S}_T, \mathbf{S}_W)$ according to (Ghojogh et al., 2019a).

3.3. Discussion on Dimensionality of the Fisher Subspace

In general, the rank of a covariance (scatter) matrix over the d -dimensional data with sample size n is at most $\min(d, n - 1)$. The d is because the covariance matrix is a $d \times d$ matrix and the n is because we iterate over n data instances for calculating the covariance matrix. The -1 is because of subtracting the mean in calculation of the

covariance matrix. For clarification, assume we only have one instance which becomes zero after removing the mean. This makes the covariance matrix a zero matrix.

According to Eq. (33), the rank of the S_W is at most $\min(d, n - 1)$ because all the instances of all the classes are considered. Hence, the rank of S_W is also at most $\min(d, n - 1)$. According to Eq. (29), the rank of the S_B is at most $\min(d, c - 1)$ because we have c iterations in its calculation.

In Eq. (57), we have $S_W^{-1}S_B$ whose rank is:

$$\begin{aligned} \text{rank}(S_W^{-1}S_B) &\leq \min(\text{rank}(S_W^{-1}), \text{rank}(S_B)) \\ &\leq \min(\min(d, n - 1), \min(d, c - 1)) \\ &= \min(d, n - 1, c - 1) \stackrel{(a)}{=} c - 1, \end{aligned} \quad (63)$$

where (a) is because we usually have $c < d, n$. Therefore, the rank of $S_W^{-1}S_B$ is limited because of the rank of S_B which is at most $c - 1$.

According to Eq. (57), the $c - 1$ leading eigenvalues will be valid and the rest are zero or very small. Therefore, the p , which is the dimensionality of the Fisher subspace, is at most $c - 1$. The $c - 1$ leading eigenvectors are considered as the Fisher directions and the rest of eigenvectors are invalid and ignored.

4. Interpretation of FDA: The Example of a Man with Weak Eyes

In this section, we interpret the FDA using a real-life example in order to better understand the essence of Fisher's method. Consider a man which has two eye problems: (1) he is color-blind and (2) his eyes are also very weak.

Suppose there are two sets of balls with red and blue colors. The man wants to discriminate the balls into red and blue classes; however, he needs help because of his eye problems.

First, consider his color-blindness. In order to help him, we separate the balls into two sets of red and blue. In other words, we increase the distances of the balls with different colors to give him a clue that which balls belong to the same class. This means that we are increasing the between-scatter of the two classes to help him.

Second, consider his very weak eyes. although the balls with different colors are almost separated, everything is blue to him. Thus, we put the balls of the same color closer to one another. In other words, we decrease the within-scatter of every class. In this way, the man sees every class as almost one blurry ball so he can discriminate the classes better.

Recall Eq. (57) which includes $S_W^{-1}S_B$. The S_B implies that we want to increase the between-scatter as we did in the first help. The S_W^{-1} implies that we want to decrease the within-scatter as done in the second help to the man.

In conclusion, FDA increases the between-scatter and decreases the within-scatter (collapses each class (Globerson & Roweis, 2006)), at the same time, for better discrimination of the classes.

5. Robust Fisher Discriminant Analysis

Robust FDA (RFDA) (Deng et al., 2007; Guo & Wang, 2015), has also addressed the problem of singularity (or close to singularity) of S_W . In RFDA, the S_W is decomposed using eigenvalue decomposition (Ghojogh et al., 2019a):

$$S_W = \Phi^T \Lambda \Phi, \quad (64)$$

where Φ and $\Lambda = \text{diag}([\lambda_1, \dots, \lambda_d]^T)$ include the eigenvectors and eigenvalues of S_W , respectively. The eigenvalues are sorted as $\lambda_1 \geq \dots \geq \lambda_d$ and the eigenvectors (columns of Φ) are sorted accordingly. If S_W is close to singularity, the first d' eigenvalues are valid and the rest $(d - d')$ eigenvalues are either very small or zero. The appropriate d' is obtained as:

$$d' := \arg \min_m \left(\frac{\sum_{j=1}^m \lambda_j}{\sum_{k=1}^d \lambda_k} \geq 0.98 \right). \quad (65)$$

In RFDA, the $(d - d')$ invalid eigenvalues are replaced with λ_* :

$$\mathbb{R}^{d \times d} \ni \Lambda' := \text{diag}([\lambda_1, \dots, \lambda_{d'}, \lambda_*, \dots, \lambda_*]^T), \quad (66)$$

where (Deng et al., 2007):

$$\lambda_* := \frac{1}{d - d'} \sum_{j=d'+1}^d \lambda_j. \quad (67)$$

Hence, the S_W is replaced with S'_W :

$$\mathbb{R}^{d \times d} \ni S'_W := \Phi^T \Lambda' \Phi, \quad (68)$$

and the robust Fisher directions are the eigenvectors of the generalized eigenvalue problem (S_B, S'_W) (Ghojogh et al., 2019a).

6. Comparison of FDA and PCA Directions

The FDA directions capture the directions where the instances of different classes fall apart and the instances in one class fall close to each other. On the other hand, the PCA directions capture the directions where the data have maximum variance (spread) regardless of the classes (Ghojogh & Crowley, 2019c). In some datasets, the FDA and PCA are orthogonal and in some datasets, they are parallel. Other cases between these two extreme cases can happen for some datasets. This depends on the spread of classes in the dataset. Figure 2 shows these cases for some two-dimensional datasets.

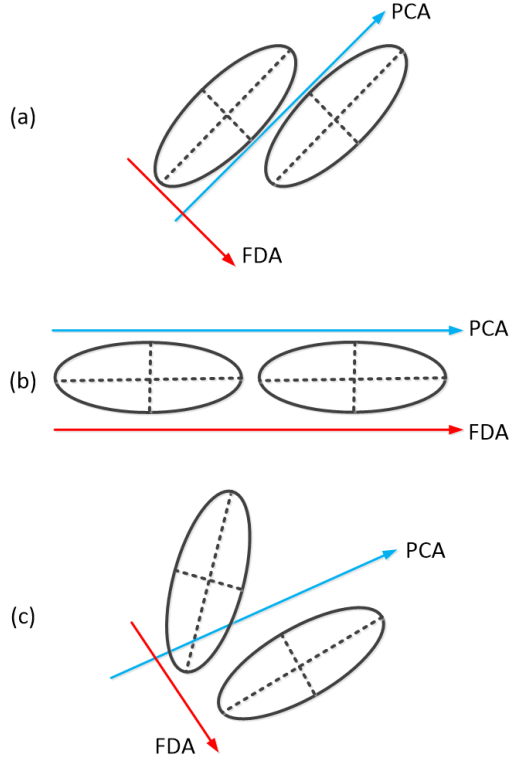


Figure 2. Comparison of FDA and PCA directions for two-dimensional data with two classes: (a) a case where FDA and PCA are orthogonal, (b) a case where FDA and PCA are equivalent (parallel), and (c) a case between the two extreme cases of (a) and (b).

Moreover, considering the Eq. (38) for S_B , the Fisher criterion becomes Eqs. (49) and (60) for one-dimensional and multi-dimensional Fisher subspaces, respectively. In these equations, the -1 is a constant and is dropped in the optimization. This has an important message about FDA: *the Fisher direction is maximizing the total variance (spread) of data, as also done in PCA, while at the same time, it minimizes the within-scatters of classes (by making use of the class labels)*. In other words, the optimization of FDA is equivalent to (we repeat Eq. (61) here):

$$\begin{aligned} & \underset{U}{\text{maximize}} \quad \text{tr}(U^\top S_T U) \\ & \text{subject to} \quad U^\top S_W U = I, \end{aligned} \quad (69)$$

while the optimization of the PCA is (Ghojogh & Crowley, 2019c):

$$\begin{aligned} & \underset{U}{\text{maximize}} \quad \text{tr}(U^\top S_T U) \\ & \text{subject to} \quad U^\top U = I. \end{aligned} \quad (70)$$

The solutions to Eqs. (69) and (70) are the generalized eigenvalue problem (S_T, S_W) and the eigenvalue problem for S_T , respectively (Ghojogh et al., 2019a).

7. FDA $\stackrel{?}{\equiv}$ LDA

The FDA is also referred to as Linear Discriminant Analysis (LDA) and Fisher LDA (FLDA). Note that FDA is a manifold (subspace) learning method and LDA (Ghojogh & Crowley, 2019a) is a classification method. However, LDA can be seen as a metric learning method (Ghojogh & Crowley, 2019a) and as metric learning is a manifold learning method (see Appendix A), there is a connection between FDA and LDA.

We know that FDA is a projection-based subspace learning method. Consider the projection vector u . According to Eq. (7), the projection of data x is:

$$x \mapsto u^\top x, \quad (71)$$

which can be done for all the data instances of every class. Thus, the mean and the covariance matrix of the class are transformed as:

$$\mu \mapsto u^\top \mu, \quad (72)$$

$$\Sigma \mapsto u^\top \Sigma u, \quad (73)$$

because of characteristics of mean and variance.

According to Eq. (42), the Fisher criterion is the ratio of the between-class variance, σ_b^2 , and within-class variance, σ_w^2 :

$$f := \frac{\sigma_b^2}{\sigma_w^2} = \frac{(u^\top \mu_2 - u^\top \mu_1)^2}{u^\top \Sigma_2 u + u^\top \Sigma_1 u} = \frac{(u^\top (\mu_2 - \mu_1))^2}{u^\top (\Sigma_2 + \Sigma_1) u}. \quad (74)$$

The FDA maximizes the Fisher criterion:

$$\underset{u}{\text{maximize}} \quad \frac{(u^\top (\mu_2 - \mu_1))^2}{u^\top (\Sigma_2 + \Sigma_1) u}, \quad (75)$$

which can be restated as:

$$\begin{aligned} & \underset{u}{\text{maximize}} \quad (u^\top (\mu_2 - \mu_1))^2, \\ & \text{subject to} \quad u^\top (\Sigma_2 + \Sigma_1) u = 1, \end{aligned} \quad (76)$$

according to Rayleigh-Ritz quotient method (Croot, 2005). The Lagrangian (Boyd & Vandenberghe, 2004) is:

$$\mathcal{L} = (u^\top (\mu_2 - \mu_1))^2 - \lambda (u^\top (\Sigma_2 + \Sigma_1) u - 1),$$

where λ is the Lagrange multiplier. Equating the derivative of \mathcal{L} to zero gives:

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial u} &= 2(\mu_2 - \mu_1)^2 u - 2\lambda(\Sigma_2 + \Sigma_1) u \stackrel{\text{set}}{=} 0 \\ \implies (\mu_2 - \mu_1)^2 u &= \lambda(\Sigma_2 + \Sigma_1) u, \end{aligned}$$

which is a generalized eigenvalue problem $((\mu_2 - \mu_1)^2, (\Sigma_2 + \Sigma_1))$ according to (Ghojogh et al., 2019a).

The projection vector is the eigenvector of $(\Sigma_2 + \Sigma_1)^{-1}(\mu_2 - \mu_1)^2$; therefore, we can say:

$$\mathbf{u} \propto (\Sigma_2 + \Sigma_1)^{-1}(\mu_2 - \mu_1)^2. \quad (77)$$

On the other hand, in LDA, the decision function is (Ghojogh & Crowley, 2019a):

$$2(\Sigma^{-1}(\mu_2 - \mu_1))^\top \mathbf{x} + (\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2) = 0, \quad (78)$$

Moreover, in LDA, the covariance matrices are assumed to be equal (Ghojogh & Crowley, 2019a): $\Sigma_1 = \Sigma_2 = \Sigma$. Therefore, in LDA, the Eq. (77) becomes (Ghojogh & Crowley, 2019a):

$$\mathbf{u} \propto (2\Sigma)^{-1}(\mu_2 - \mu_1)^2 \propto \Sigma^{-1}(\mu_2 - \mu_1)^2. \quad (79)$$

According to Eq. (71), we have:

$$\mathbf{u}^\top \mathbf{x} \propto (\Sigma^{-1}(\mu_2 - \mu_1)^2)^\top \mathbf{x}. \quad (80)$$

Comparing Eq. (80) with Eq. (78) shows that LDA and FDA are equivalent up to a scaling factor $(\mu_1 - \mu_2)^\top \Sigma^{-1}(\mu_1 - \mu_2)$. Note that this term is multiplied as an exponential factor before taking logarithm to obtain Eq. (78), so this term a scaling factor (see (Ghojogh & Crowley, 2019a) for more details). It should be noted that in manifold (subspace) learning, the scale does not matter because all the distances scale similarly. Hence, we can say that LDA and FDA are equivalent:

$$\text{LDA} \equiv \text{FDA}. \quad (81)$$

Therefore, *the two subspaces of FDA and LDA are the same subspace*. This sheds light to why LDA and FDA are used interchangeably in the literature.

Note that LDA assumes *one* (and not several) Gaussian for every class (Ghojogh & Crowley, 2019a) and so does the FDA because they are equivalent. That is why FDA faces problem for multi-modal data (Sugiyama, 2007).

8. Fisher Forest

If the data include several different types of data which may even have different dimensionality. Some examples of these types of data are different key-poses in action human action recognition or different facial expressions that a face can have. In this case, we can use the *Fisher forest* (Ghojogh & Mohammadzade, 2017). Note that forest here does not imply an ensemble of trees but means an ensemble of the Fisher subspaces.

Let the number of the data types be z and let the dimensionality of the ℓ -th data type be $d_{|\ell}$. We usually have a dataset $\{\mathbf{x}_i\}_{i=1}^n$ where $\mathbf{x}_i \in \mathbb{R}^d$. Every type of data is the whole dataset but having only a subset of the features, i.e.,

$\{\mathbf{x}_{i|\ell}\}_{i=1}^n$ where $\mathbf{x}_{i|\ell} \in \mathbb{R}^{d_{|\ell}}$. The features of the ℓ -th type are a subset of the features of the dataset, i.e., $d_{|\ell} \leq d$. Note that, we do not necessarily have the same value for $d_{|\ell}$ in all the data types. The i -th instance of the j -th class having the ℓ -th type is denoted by $\mathbf{x}_{i|\ell}^{(j)}$.

For example, in the key-poses of human action, the important of skeletal joints can be different in various key-poses (Ghojogh & Mohammadzade, 2017). Thus, some joints are taken in a specific key-joint and some other are taken in another key-pose. Note that a key-pose can have five key-joints but another key-pose can have three key-joints. Another example is using different landmarks for different facial expressions; for example, eye-brows, lips, and chin for wondering but just lips for smiling. As can be seen, Fisher forest can be useful for handling the data types with different features and even dimensionality.

The between- and within-scatters for the ℓ -th data types (for all $\ell \in \{1, \dots, z\}$) are defined as (Ghojogh & Mohammadzade, 2017):

$$\mathbb{R}^{d_\ell \times d_\ell} \ni \mathbf{S}_{B|\ell} := \sum_{j=1}^c n_j (\mu_{j|\ell} - \mu_{|\ell})(\mu_{j|\ell} - \mu_{|\ell})^\top, \quad (82)$$

$$\mathbb{R}^{d_\ell \times d_\ell} \ni \mathbf{S}_{W|\ell} := \sum_{j=1}^c \sum_{i=1}^{n_j} (\mathbf{x}_{i|\ell}^{(j)} - \mu_{j|\ell})(\mathbf{x}_{i|\ell}^{(j)} - \mu_{j|\ell})^\top, \quad (83)$$

where:

$$\mathbb{R}^{d_{|\ell}} \ni \mu_{j|\ell} := \frac{1}{n_j} \sum_{i=1}^{n_j} \mathbf{x}_{i|\ell}^{(j)}, \quad (84)$$

$$\mathbb{R}^{d_{|\ell}} \ni \mu_{|\ell} := \frac{1}{n} \sum_{i=1}^n \mathbf{x}_{i|\ell}. \quad (85)$$

The ℓ -th Fisher subspace is spanned by the eigenvectors of $\mathbf{S}_{W|\ell}^{-1} \mathbf{S}_{B|\ell}$.

Hence, z Fisher subspaces are trained. In the test phase, the data instance is projected onto every subspace. If we want to classify the data instance in the projected subspaces, we will have z classification results after projection onto these z subspaces. We can then use majority voting for a final classification of the data instance (Ghojogh & Mohammadzade, 2017). The effectiveness of the majority voting can be explained because of ensemble learning (Polikar, 2012; Ghojogh & Crowley, 2019b). We can also normalize the distances in the subspaces of Fisher forest for the sake of classification (see (Ghojogh & Mohammadzade, 2017) for more details).

9. Kernel Fisher Discriminant Analysis

9.1. Kernels and Hilbert Space

Suppose that $\phi : \mathcal{X} \rightarrow \mathcal{H}$ is a function which maps the data \mathbf{x} to Hilbert space (feature space). The ϕ is called *pulling function*. In other words, $\mathbf{x} \mapsto \phi(\mathbf{x})$. Let t denote the dimensionality of the feature space, i.e., $\phi(\mathbf{x}) \in \mathbb{R}^t$ while $\mathbf{x} \in \mathbb{R}^d$. Note that we usually have $t \gg d$.

If \mathcal{X} denotes the set of points, i.e., $\mathbf{x} \in \mathcal{X}$, the kernel of two vectors \mathbf{x}_1 and \mathbf{x}_2 is $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and is defined as (Hofmann et al., 2008; Herbrich, 2001):

$$k(\mathbf{x}_1, \mathbf{x}_2) := \phi(\mathbf{x}_1)^\top \phi(\mathbf{x}_2), \quad (86)$$

which is a measure of *similarity* between the two vectors because the inner product captures similarity.

We can compute the kernel of two matrices $\mathbf{X}_1 \in \mathbb{R}^{d \times n_1}$ and $\mathbf{X}_2 \in \mathbb{R}^{d \times n_2}$ and have a *kernel matrix* (also called *Gram matrix*):

$$\mathbb{R}^{n_1 \times n_2} \ni \mathbf{K}(\mathbf{X}_1, \mathbf{X}_2) := \Phi(\mathbf{X}_1)^\top \Phi(\mathbf{X}_2), \quad (87)$$

where $\Phi(\mathbf{X}_1) := [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_{n_1})] \in \mathbb{R}^{t \times n_1}$ is the matrix of mapped \mathbf{X}_1 to the feature space. The $\Phi(\mathbf{X}_2) \in \mathbb{R}^{t \times n_2}$ is defined similarly. We can compute the kernel matrix of dataset $\mathbf{X} \in \mathbb{R}^{d \times n}$ over itself:

$$\mathbb{R}^{n \times n} \ni \mathbf{K}_x := \mathbf{K}(\mathbf{X}, \mathbf{X}) = \Phi(\mathbf{X})^\top \Phi(\mathbf{X}), \quad (88)$$

where $\Phi(\mathbf{X}) := [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)] \in \mathbb{R}^{t \times n}$ is the pulled (mapped) data.

Note that in kernel methods, the pulled data $\Phi(\mathbf{X})$ are usually not available and merely the kernel matrix $\mathbf{K}(\mathbf{X}, \mathbf{X})$, which is the inner product of the pulled data with itself, is available.

There exist different types of kernels. Some of the most well-known kernels are:

$$\text{Linear: } k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^\top \mathbf{x}_2 + c_1, \quad (89)$$

$$\text{Polynomial: } k(\mathbf{x}_1, \mathbf{x}_2) = (c_1 \mathbf{x}_1^\top \mathbf{x}_2 + c_2)^{c_3}, \quad (90)$$

$$\text{Gaussian: } k(\mathbf{x}_1, \mathbf{x}_2) = \exp\left(-\frac{\|\mathbf{x}_1 - \mathbf{x}_2\|_2^2}{2\sigma^2}\right), \quad (91)$$

$$\text{Sigmoid: } k(\mathbf{x}_1, \mathbf{x}_2) = \tanh(c_1 \mathbf{x}_1^\top \mathbf{x}_2 + c_2), \quad (92)$$

where c_1, c_2, c_3 , and σ are scalar constants. The Gaussian and Sigmoid kernels are also called Radial Basis Function (RBF) and hyperbolic tangent, respectively. Note that the Gaussian kernel can also be written as $\exp(-\gamma \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2)$ where $\gamma > 0$.

It is noteworthy to mention that in the RBF kernel, the dimensionality of the feature space is infinite. The reason lies in the Maclaurin series expansion (Taylor series expansion at zero) of this kernel:

$$\exp(-\gamma r) \approx 1 - \gamma r + \frac{\gamma^2}{2!} r^2 - \frac{\gamma^3}{3!} r^3 + \dots,$$

where $r := \|\mathbf{x}_1 - \mathbf{x}_2\|_2^2$, which is infinite dimensional with respect to r .

9.2. One-dimensional Subspace

9.2.1. SCATTERS IN TWO-CLASS CASE

The Eq. (26) in the feature space is:

$$\mathbb{R}^{t \times t} \ni \Phi(\mathbf{S}_B) := (\phi(\boldsymbol{\mu}_1) - \phi(\boldsymbol{\mu}_2))(\phi(\boldsymbol{\mu}_1) - \phi(\boldsymbol{\mu}_2))^\top, \quad (93)$$

where the mean of the j -th class in the feature space is:

$$\mathbb{R}^t \ni \phi(\boldsymbol{\mu}_j) := \frac{1}{n_j} \sum_{i=1}^{n_j} \phi(\mathbf{x}_i^{(j)}). \quad (94)$$

According to the representation theory (Alperin, 1993), any solution (direction) $\phi(\mathbf{u}) \in \mathcal{H}$ must lie in the span of “all” the training vectors mapped to \mathcal{H} , i.e., $\Phi(\mathbf{X}) = [\phi(\mathbf{x}_1), \dots, \phi(\mathbf{x}_n)] \in \mathbb{R}^{t \times n}$ (usually $t \gg d$). Note that \mathcal{H} denotes the Hilbert space (feature space). Therefore, we can state that:

$$\mathbb{R}^t \ni \phi(\mathbf{u}) = \sum_{i=1}^n \theta_i \phi(\mathbf{x}_i) = \Phi(\mathbf{X}) \boldsymbol{\theta}, \quad (95)$$

where $\mathbb{R}^n \ni \boldsymbol{\theta} := [\theta_1, \dots, \theta_n]^\top$ is the unknown vector of coefficients, and $\phi(\mathbf{u}) \in \mathbb{R}^t$ is the pulled Fisher direction to the feature space. The pulled directions can be put together in $\mathbb{R}^{t \times p} \ni \Phi(\mathbf{U}) := [\phi(\mathbf{u}_1), \dots, \phi(\mathbf{u}_p)]$:

$$\mathbb{R}^{t \times p} \ni \Phi(\mathbf{U}) = \Phi(\mathbf{X}) \boldsymbol{\Theta}, \quad (96)$$

where $\boldsymbol{\Theta} := [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p] \in \mathbb{R}^{n \times p}$.

The d_B in the feature space is:

$$\mathbb{R} \ni d_B := \phi(\mathbf{u})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{u}) \quad (97)$$

$$\stackrel{(a)}{=} \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top (\phi(\boldsymbol{\mu}_1) - \phi(\boldsymbol{\mu}_2)) (\phi(\boldsymbol{\mu}_1) - \phi(\boldsymbol{\mu}_2))^\top \Phi(\mathbf{X}) \boldsymbol{\theta}, \quad (98)$$

where (a) is because of Eqs. (93). and (95).

For the j -th class (here $j \in \{1, 2\}$), we have:

$$\begin{aligned} \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \phi(\boldsymbol{\mu}_j) &\stackrel{(95)}{=} \sum_{i=1}^n \theta_i \phi(\mathbf{x}_i)^\top \phi(\boldsymbol{\mu}_j) \\ &\stackrel{(94)}{=} \frac{1}{n_j} \sum_{i=1}^n \sum_{k=1}^{n_j} \theta_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_k^{(j)}) \\ &\stackrel{(86)}{=} \frac{1}{n_j} \sum_{i=1}^n \sum_{k=1}^{n_j} \theta_i k(\mathbf{x}_i, \mathbf{x}_k^{(j)}) = \boldsymbol{\theta}^\top \mathbf{m}_j, \end{aligned} \quad (99)$$

where $\mathbf{m}_j \in \mathbb{R}^n$ whose i -th entry is:

$$\mathbf{m}_j(i) := \frac{1}{n_j} \sum_{k=1}^{n_j} k(\mathbf{x}_i, \mathbf{x}_k^{(j)}). \quad (100)$$

Hence, Eq. (98) becomes:

$$d_B \stackrel{(99)}{=} \boldsymbol{\theta}^\top (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top \boldsymbol{\theta} = \boldsymbol{\theta}^\top \mathbf{M} \boldsymbol{\theta}, \quad (101)$$

where:

$$\mathbb{R}^{n \times n} \ni \mathbf{M} := (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top, \quad (102)$$

is the *between-scatter* in kernel FDA. Hence, the Eq. (98) becomes:

$$d_B = \phi(\mathbf{u})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{u}) = \boldsymbol{\theta}^\top \mathbf{M} \boldsymbol{\theta}. \quad (103)$$

The Eq. (33) in the feature space is:

$$\begin{aligned} \mathbb{R}^{t \times t} \ni \Phi(\mathbf{S}_W) &:= \\ &\sum_{j=1}^c \sum_{i=1}^{n_j} (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j)) (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j))^\top. \end{aligned} \quad (104)$$

The d_W in the feature space is:

$$\begin{aligned} \mathbb{R} \ni d_W &:= \phi(\mathbf{u})^\top \Phi(\mathbf{S}_W) \phi(\mathbf{u}) \\ &\stackrel{(a)}{=} \left(\sum_{\ell=1}^n \theta_\ell \phi(\mathbf{x}_\ell)^\top \right) \left(\sum_{j=1}^c \sum_{i=1}^{n_j} (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j)) \right. \\ &\quad \left. (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j))^\top \right) \left(\sum_{k=1}^n \theta_k \phi(\mathbf{x}_k) \right) \\ &= \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \left(\theta_\ell \phi(\mathbf{x}_\ell)^\top (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j)) \right. \\ &\quad \left. (\phi(\mathbf{x}_i^{(j)}) - \phi(\boldsymbol{\mu}_j))^\top \theta_k \phi(\mathbf{x}_k) \right) \\ &\stackrel{(94)}{=} \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \\ &\quad \left(\theta_\ell \phi(\mathbf{x}_\ell)^\top (\phi(\mathbf{x}_i^{(j)}) - \frac{1}{n_j} \sum_{e=1}^{n_j} \phi(\mathbf{x}_e^{(j)})) \right. \\ &\quad \left. (\phi(\mathbf{x}_i^{(j)}) - \frac{1}{n_j} \sum_{z=1}^{n_j} \phi(\mathbf{x}_z^{(j)}))^\top \theta_k \phi(\mathbf{x}_k) \right) \\ &\stackrel{(86)}{=} \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \\ &\quad \left(\theta_\ell k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) - \frac{1}{n_j} \sum_{e=1}^{n_j} \theta_\ell k(\mathbf{x}_\ell, \mathbf{x}_e^{(j)}) \right) \\ &\quad \left(\theta_k k(\mathbf{x}_i^{(j)}, \mathbf{x}_k) - \frac{1}{n_j} \sum_{z=1}^{n_j} \theta_k k(\mathbf{x}_z^{(j)}, \mathbf{x}_k) \right) \end{aligned}$$

$$\begin{aligned} &\stackrel{(b)}{=} \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \\ &\quad \left(\theta_\ell k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) - \frac{1}{n_j} \sum_{e=1}^{n_j} \theta_\ell k(\mathbf{x}_\ell, \mathbf{x}_e^{(j)}) \right) \\ &\quad \left(\theta_k k(\mathbf{x}_k, \mathbf{x}_i^{(j)}) - \frac{1}{n_j} \sum_{z=1}^{n_j} \theta_k k(\mathbf{x}_k, \mathbf{x}_z^{(j)}) \right) \\ &= \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \\ &\quad \left(\theta_\ell \theta_k k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_i^{(j)}) \right. \\ &\quad \left. - \frac{2 \theta_\ell \theta_k}{n_j} \sum_{z=1}^{n_j} k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_z^{(j)}) \right. \\ &\quad \left. + \frac{\theta_\ell \theta_k}{n_j^2} \sum_{e=1}^{n_j} \sum_{z=1}^{n_j} k(\mathbf{x}_\ell, \mathbf{x}_e^{(j)}) k(\mathbf{x}_k, \mathbf{x}_z^{(j)}) \right) \\ &= \sum_{j=1}^c \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \\ &\quad \left(\theta_\ell \theta_k k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_i^{(j)}) \right. \\ &\quad \left. - \frac{\theta_\ell \theta_k}{n_j} \sum_{z=1}^{n_j} k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_z^{(j)}) \right) \\ &= \sum_{j=1}^c \left(\sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \left(\theta_\ell \theta_k k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_i^{(j)}) \right) \right. \\ &\quad \left. - \sum_{\ell=1}^n \sum_{i=1}^{n_j} \sum_{k=1}^n \left(\frac{\theta_\ell \theta_k}{n_j} \sum_{z=1}^{n_j} k(\mathbf{x}_\ell, \mathbf{x}_i^{(j)}) k(\mathbf{x}_k, \mathbf{x}_z^{(j)}) \right) \right) \\ &\stackrel{(c)}{=} \sum_{j=1}^c (\boldsymbol{\theta}^\top \mathbf{K}_j \mathbf{K}_j^\top \boldsymbol{\theta} - \boldsymbol{\theta}^\top \mathbf{K}_j \frac{1}{n_j} \mathbf{1} \mathbf{1}^\top \mathbf{K}_j^\top \boldsymbol{\theta}) \\ &= \sum_{j=1}^c \boldsymbol{\theta}^\top \mathbf{K}_j \left(\mathbf{I} - \frac{1}{n_j} \mathbf{1} \mathbf{1}^\top \right) \mathbf{K}_j^\top \boldsymbol{\theta} \\ &\stackrel{(d)}{=} \sum_{j=1}^c \boldsymbol{\theta}^\top \mathbf{K}_j \mathbf{H}_j \mathbf{K}_j^\top \boldsymbol{\theta} = \boldsymbol{\theta}^\top \left(\sum_{j=1}^c \mathbf{K}_j \mathbf{H}_j \mathbf{K}_j^\top \right) \boldsymbol{\theta}, \end{aligned}$$

where (a) is because of Eqs. (104) and (95), (b) is because $k(\mathbf{x}_1, \mathbf{x}_2) = k(\mathbf{x}_2, \mathbf{x}_1) \in \mathbb{R}$, and (c) is because $\mathbf{K}_j \in \mathbb{R}^{n \times n_j}$ is the kernel matrix of the whole training data and the training data of the j -th class. The (a, b)-th element of \mathbf{K}_j is:

$$\mathbf{K}_j(a, b) := k(\mathbf{x}_a, \mathbf{x}_b^{(j)}). \quad (105)$$

The (d) is because:

$$\mathbb{R}^{n_j \times n_j} \ni \mathbf{H}_j := \mathbf{I} - \frac{1}{n_j} \mathbf{1}\mathbf{1}^\top, \quad (106)$$

is the *centering matrix* (see Appendix A in (Ghojogh & Crowley, 2019c)).

We define:

$$\mathbb{R}^{n \times n} \ni \mathbf{N} := \sum_{j=1}^c \mathbf{K}_j \mathbf{H}_j \mathbf{K}_j^\top, \quad (107)$$

as the *within-scatter* in kernel FDA. Hence, the d_W becomes:

$$d_W = \phi(\mathbf{u})^\top \Phi(\mathbf{S}_W) \phi(\mathbf{u}) = \boldsymbol{\theta}^\top \mathbf{N} \boldsymbol{\theta}. \quad (108)$$

The kernel Fisher criterion is:

$$f(\boldsymbol{\theta}) := \frac{d_B(\boldsymbol{\theta})}{d_W(\boldsymbol{\theta})} = \frac{\phi(\mathbf{u})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{u})}{\phi(\mathbf{u})^\top \Phi(\mathbf{S}_W) \phi(\mathbf{u})} = \frac{\boldsymbol{\theta}^\top \mathbf{M} \boldsymbol{\theta}}{\boldsymbol{\theta}^\top \mathbf{N} \boldsymbol{\theta}}, \quad (109)$$

where the $\boldsymbol{\theta} \in \mathbb{R}^n$ is the *kernel Fisher direction*.

Similar to the solution of Eq. (42), the solution to maximization of Eq. (109) is:

$$\mathbf{M} \boldsymbol{\theta} = \lambda \mathbf{N} \boldsymbol{\theta}, \quad (110)$$

which is a generalized eigenvalue problem (\mathbf{M}, \mathbf{N}) according to (Ghojogh et al., 2019a). The $\boldsymbol{\theta}$ is the eigenvector with the largest eigenvalue (because the optimization is maximization) and the λ is the corresponding eigenvalue. The $\boldsymbol{\theta}$ is the *kernel Fisher direction* or *kernel Fisher axis*.

Again, one possible solution to the generalized eigenvalue problem (\mathbf{M}, \mathbf{N}) is (Ghojogh et al., 2019a):

$$\boldsymbol{\theta} = \mathbf{eig}(\mathbf{N}^{-1} \mathbf{M}), \quad (111)$$

or (Ghojogh et al., 2019a):

$$\boldsymbol{\theta} = \mathbf{eig}((\mathbf{N} + \varepsilon \mathbf{I})^{-1} \mathbf{M}), \quad (112)$$

where $\mathbf{eig}(\cdot)$ denotes the eigenvector of the matrix with the largest eigenvalue.

The projection and reconstruction of the training data point \mathbf{x}_i and the out-of-sample data point \mathbf{x}_t are:

$$\begin{aligned} \mathbb{R} \ni \phi(\tilde{\mathbf{x}}_i) &= \phi(\mathbf{u})^\top \phi(\mathbf{x}_i) \stackrel{(95)}{=} \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \phi(\mathbf{x}_i) \\ &= \boldsymbol{\theta}^\top \mathbf{k}(\mathbf{X}, \mathbf{x}_i), \end{aligned} \quad (113)$$

$$\begin{aligned} \mathbb{R}^t \ni \phi(\hat{\mathbf{x}}_i) &= \phi(\mathbf{u}) \phi(\mathbf{u})^\top \phi(\mathbf{x}_i) \\ &\stackrel{(95)}{=} \Phi(\mathbf{X}) \boldsymbol{\theta} \boldsymbol{\theta}^\top \mathbf{k}(\mathbf{X}, \mathbf{x}_i), \end{aligned} \quad (114)$$

$$\mathbb{R} \ni \phi(\tilde{\mathbf{x}}_t) = \boldsymbol{\theta}^\top \mathbf{k}(\mathbf{X}, \mathbf{x}_t), \quad (115)$$

$$\mathbb{R}^t \ni \phi(\hat{\mathbf{x}}_t) = \Phi(\mathbf{X}) \boldsymbol{\theta} \boldsymbol{\theta}^\top \mathbf{k}(\mathbf{X}, \mathbf{x}_t). \quad (116)$$

However, in reconstruction expressions, the $\Phi(\mathbf{X})$ is not necessarily available; therefore, in kernel FDA, similar to kernel PCA (Ghojogh & Crowley, 2019c), *reconstruction cannot be done*. For the whole training and out-of-sample data, the projections are:

$$\mathbb{R}^{1 \times n} \ni \Phi(\tilde{\mathbf{X}}) = \boldsymbol{\theta}^\top \mathbf{K}(\mathbf{X}, \mathbf{X}), \quad (117)$$

$$\mathbb{R}^{1 \times n_t} \ni \Phi(\tilde{\mathbf{X}}_t) = \boldsymbol{\theta}^\top \mathbf{K}(\mathbf{X}, \mathbf{X}_t). \quad (118)$$

9.2.2. SCATTERS IN MULTI-CLASS CASE: VARIANT 1

In multi-class case for kernel FDA, the *within-scatter* is the same as in the two-class case, which is Eq. (107) and d_W is also Eq. (108). However, the between-scatter is different. The between-scatter, Eq. (29), in the feature space is:

$$\begin{aligned} \mathbb{R}^{t \times t} \ni \Phi(\mathbf{S}_B) &:= \\ &\sum_{j=1}^c (\phi(\boldsymbol{\mu}_j) - \phi(\boldsymbol{\mu})) (\phi(\boldsymbol{\mu}_j) - \phi(\boldsymbol{\mu}))^\top, \end{aligned} \quad (119)$$

where the total mean in the feature space is:

$$\mathbb{R}^t \ni \phi(\boldsymbol{\mu}) := \frac{1}{\sum_{k=1}^c n_k} \sum_{j=1}^c n_j \phi(\boldsymbol{\mu}_j) = \frac{1}{n} \sum_{i=1}^n \phi(\mathbf{x}_i), \quad (120)$$

The d_B in the feature space is:

$$\begin{aligned} \mathbb{R} \ni d_B &:= \phi(\mathbf{u})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{u}) \\ &\stackrel{(a)}{=} \sum_{j=1}^c \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top (\phi(\boldsymbol{\mu}_j) - \phi(\boldsymbol{\mu})) \\ &\quad (\phi(\boldsymbol{\mu}_j) - \phi(\boldsymbol{\mu}))^\top \Phi(\mathbf{X}) \boldsymbol{\theta}, \end{aligned} \quad (121)$$

where (a) is because of Eqs. (119) and (95). We have:

$$\begin{aligned} \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \phi(\boldsymbol{\mu}) &\stackrel{(95)}{=} \sum_{i=1}^n \theta_i \phi(\mathbf{x}_i)^\top \phi(\boldsymbol{\mu}) \\ &\stackrel{(120)}{=} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^n \theta_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_k) \\ &\stackrel{(86)}{=} \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^n \theta_i k(\mathbf{x}_i, \mathbf{x}_k) = \boldsymbol{\theta}^\top \mathbf{m}_*, \end{aligned} \quad (122)$$

where $\mathbf{m}_* \in \mathbb{R}^n$ whose i -th entry is:

$$\mathbf{m}_*(i) := \frac{1}{n} \sum_{k=1}^n k(\mathbf{x}_i, \mathbf{x}_k). \quad (123)$$

According to Eqs. (99) and (122), the Eq. (121) becomes:

$$d_B = \boldsymbol{\theta}^\top \sum_{j=1}^c (\mathbf{m}_j - \mathbf{m}_*)(\mathbf{m}_j - \mathbf{m}_*)^\top \boldsymbol{\theta} = \boldsymbol{\theta}^\top \mathbf{M} \boldsymbol{\theta}, \quad (124)$$

where:

$$\mathbb{R}^{n \times n} \ni \mathbf{M} := \sum_{j=1}^c (\mathbf{m}_j - \mathbf{m}_*) (\mathbf{m}_j - \mathbf{m}_*)^\top, \quad (125)$$

is the *between-scatter* in kernel FDA. Similar to Eq. (31), some researches consider the following instead:

$$\mathbb{R}^{n \times n} \ni \mathbf{M} := \sum_{j=1}^c n_j (\mathbf{m}_j - \mathbf{m}_*) (\mathbf{m}_j - \mathbf{m}_*)^\top. \quad (126)$$

Hence, the Eq. (121) becomes:

$$d_B = \phi(\mathbf{u})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{u}) = \boldsymbol{\theta}^\top \mathbf{M} \boldsymbol{\theta}, \quad (127)$$

where \mathbf{M} here is Eq. (125) or (126).

The Fisher direction is again Eq. (109) and the solution is again the generalized eigenvalue problem (\mathbf{M}, \mathbf{N}) according to (Ghojogh et al., 2019a).

9.2.3. SCATTERS IN MULTI-CLASS CASE: VARIANT 2

Again, in the second version of multi-class case for kernel FDA, the *within-scatter* is the same as in the two-class case, which is Eq. (107) and d_W is also Eq. (108).

For the between scatter in the second version, we start with the Eqs. (49) and (50). We kernelize the objective function of the Eq. (50):

$$d_T := \phi(\mathbf{u})^\top \Phi(\mathbf{S}_T) \phi(\mathbf{u}), \quad (128)$$

where total-scatter, Eq. (36), is pulled as:

$$\begin{aligned} \mathbb{R}^{t \times t} \ni \Phi(\mathbf{S}_T) &:= \\ &\sum_{k=1}^n (\phi(\mathbf{x}_k) - \phi(\boldsymbol{\mu})) (\phi(\mathbf{x}_k) - \phi(\boldsymbol{\mu}))^\top. \end{aligned} \quad (129)$$

According to Eqs. (95), (128), and (129), we have:

$$\begin{aligned} d_T &= \sum_{k=1}^n \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top (\phi(\mathbf{x}_k) - \phi(\boldsymbol{\mu})) \\ &\quad (\phi(\mathbf{x}_k) - \phi(\boldsymbol{\mu}))^\top \Phi(\mathbf{X}) \boldsymbol{\theta}. \end{aligned}$$

According to Eq. (122), we have:

$$\boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \phi(\boldsymbol{\mu}) = \boldsymbol{\theta}^\top \mathbf{m}_*, \quad (130)$$

where \mathbf{m}_* is Eq. (123). On the other hand, we have:

$$\begin{aligned} \boldsymbol{\theta}^\top \Phi(\mathbf{X})^\top \phi(\mathbf{x}_k) &\stackrel{(95)}{=} \sum_{i=1}^n \theta_i \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_k) \\ &\stackrel{(86)}{=} \sum_{i=1}^n \theta_i k(\mathbf{x}_i, \mathbf{x}_k) = \boldsymbol{\theta}^\top \mathbf{g}_k, \end{aligned} \quad (131)$$

where $\mathbf{g}_k \in \mathbb{R}^n$ whose i -th entry is:

$$\mathbf{g}_k(i) := k(\mathbf{x}_i, \mathbf{x}_k). \quad (132)$$

Hence:

$$d_T = \sum_{k=1}^n \boldsymbol{\theta}^\top (\mathbf{g}_k - \mathbf{m}_*) (\mathbf{g}_k - \mathbf{m}_*)^\top \boldsymbol{\theta} = \boldsymbol{\theta}^\top \mathbf{G} \boldsymbol{\theta}, \quad (133)$$

where:

$$\mathbb{R}^{n \times n} \ni \mathbf{G} := \sum_{k=1}^n (\mathbf{g}_k - \mathbf{m}_*) (\mathbf{g}_k - \mathbf{m}_*)^\top. \quad (134)$$

The denominator of the Fisher criterion in the feature space is again the Eq. (108).

The optimization will be similar to Eq. (50) but in the feature space:

$$\begin{aligned} &\underset{\boldsymbol{\theta}}{\text{maximize}} \quad \boldsymbol{\theta}^\top \mathbf{G} \boldsymbol{\theta} \\ &\text{subject to} \quad \boldsymbol{\theta}^\top \mathbf{N} \boldsymbol{\theta} = 1, \end{aligned} \quad (135)$$

whose solution is similarly obtained as:

$$\mathbf{G} \boldsymbol{\theta} = \lambda \mathbf{N} \boldsymbol{\theta}, \quad (136)$$

which is a generalized eigenvalue problem (\mathbf{G}, \mathbf{N}) according to (Ghojogh et al., 2019a).

9.3. Multi-dimensional Subspace

In the previous section, the one-dimensional kernel Fisher subspace was discussed. In multi-dimensional kernel Fisher subspace, the within- and between-scatters are the same but the fisher criterion is different. According to Eq. (96), the d_B and d_W are:

$$d_B = \text{tr}(\phi(\mathbf{U})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{U})) = \text{tr}(\boldsymbol{\Theta}^\top \mathbf{M} \boldsymbol{\Theta}), \quad (137)$$

$$d_W = \text{tr}(\phi(\mathbf{U})^\top \Phi(\mathbf{S}_W) \phi(\mathbf{U})) = \text{tr}(\boldsymbol{\Theta}^\top \mathbf{N} \boldsymbol{\Theta}), \quad (138)$$

where $\mathbb{R}^{n \times p} \ni \boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_p]$ and $\mathbf{M} \in \mathbb{R}^{n \times n}$ and $\mathbf{N} \in \mathbb{R}^{n \times n}$ are the between- and within-scatters, respectively, determined for either two-class or multi-class case.

The Fisher criterion becomes:

$$\begin{aligned} f(\boldsymbol{\Theta}) &:= \frac{d_B(\boldsymbol{\Theta})}{d_W(\boldsymbol{\Theta})} = \frac{\text{tr}(\phi(\mathbf{U})^\top \Phi(\mathbf{S}_B) \phi(\mathbf{U}))}{\text{tr}(\phi(\mathbf{U})^\top \Phi(\mathbf{S}_W) \phi(\mathbf{U}))} \\ &= \frac{\text{tr}(\boldsymbol{\Theta}^\top \mathbf{M} \boldsymbol{\Theta})}{\text{tr}(\boldsymbol{\Theta}^\top \mathbf{N} \boldsymbol{\Theta})}, \end{aligned} \quad (139)$$

where the columns of $\boldsymbol{\Theta}$ are the *kernel Fisher directions*. Similar to Eq. (54), the solution to maximization of this criterion is:

$$\mathbf{M} \boldsymbol{\Theta} = \mathbf{N} \boldsymbol{\Theta} \boldsymbol{\Lambda}, \quad (140)$$

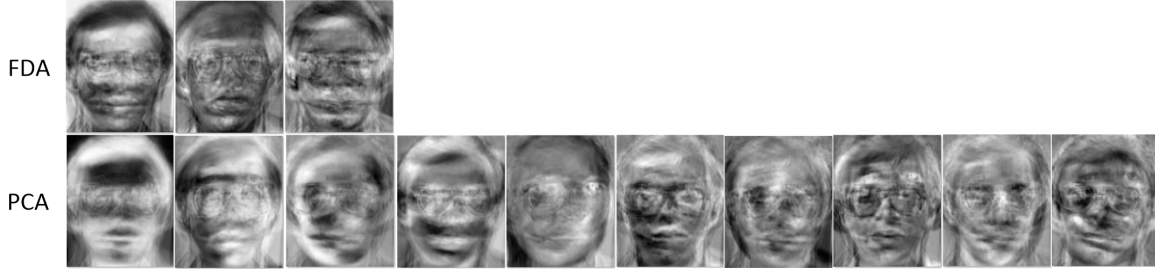


Figure 3. The projection directions (ghost faces) of FDA and PCA for the first four classes of facial AT&T dataset.

the generalized eigenvalue problem (M, N) according to (Ghojogh et al., 2019a). The columns of Θ are the eigenvectors sorted from the largest to smallest eigenvalues (because the optimization is maximization) and the diagonal entries of Λ are the corresponding eigenvalues.

Again, we can have another variant of kernel FDA for the multi-dimensional sub-space where the optimization is (similar to Eq. (135)):

$$\begin{aligned} & \underset{\Theta}{\text{maximize}} \quad \text{tr}(\Theta^\top G \Theta) \\ & \text{subject to} \quad \Theta^\top N \Theta = I, \end{aligned} \quad (141)$$

whose solution is similarly obtained as:

$$G\Theta = N\Theta\Lambda, \quad (142)$$

which is a generalized eigenvalue problem (G, N) according to (Ghojogh et al., 2019a).

As mentioned before, in kernel FDA, we do not have reconstruction. The projection of the training data point x_i and the out-of-sample data point x_t are:

$$\begin{aligned} \mathbb{R}^p \ni \phi(\tilde{x}_i) &= \Phi(U)^\top \phi(x_i) \stackrel{(96)}{=} \Theta^\top \Phi(X)^\top \phi(x_i) \\ &= \Theta^\top k(X, x_i), \end{aligned} \quad (143)$$

$$\mathbb{R}^p \ni \phi(\tilde{x}_t) = \Theta^\top k(X, x_t). \quad (144)$$

For the whole training and out-of-sample data, the projections are:

$$\mathbb{R}^{p \times n} \ni \Phi(\tilde{X}) = \Theta^\top K(X, X), \quad (145)$$

$$\mathbb{R}^{p \times n_t} \ni \Phi(\tilde{X}_t) = \Theta^\top K(X, X_t). \quad (146)$$

9.4. Discussion on Dimensionality of the Kernel Fisher Subspace

According to Eq. (107), the rank of the N is at most $\min(n, c)$ because the matrix is $n \times n$ and its calculation includes c iterations. Hence, the rank of N^{-1} is also at most $\min(n, c)$. According to Eq. (126), the rank of the M is at most $\min(n, c - 1)$ because the matrix is $n \times n$, we have c iterations in its calculation, and -1 is because

of subtracting the mean (refer to the explanation in Section 3.3).

In Eq. (111), we have $N^{-1}M$ whose rank is:

$$\begin{aligned} \text{rank}(N^{-1}M) &\leq \min(\text{rank}(N^{-1}), \text{rank}(M)) \\ &\leq \min(\min(n, c), \min(n, c - 1)) \\ &= \min(n, c - 1) \stackrel{(a)}{=} c - 1, \end{aligned} \quad (147)$$

where (a) is because we usually have $c < n$. Therefore, the rank of $N^{-1}M$ is limited because of the rank of M which is at most $c - 1$.

According to Eq. (111), the $c - 1$ leading eigenvalues will be valid and the rest are zero or very small. Therefore, the p , which is the dimensionality of the kernel Fisher subspace, is at most $c - 1$. The $c - 1$ leading eigenvectors are considered as the kernel Fisher directions and the rest of eigenvectors are invalid and ignored.

10. Simulations

For the simulations, we used the AT&T face dataset which includes 400 images, 40 subjects, and 10 images per subject. The images of every person have different poses and expressions. For better visualization of separation of classes in the projection subspace, we only used the images of the first four subjects. The dataset, except for reconstruction experiments, was standardized so that its mean and variance became zero and one, respectively.

10.1. Visualization of the Projection Directions

First, we used the entire 40 images for training FDA, kernel FDA, PCA, and kernel PCA where the used kernels are linear, Radial Basis Function (RBF), and cosine kernels. As we have four classes, the number of FDA directions is three. The three FDA directions and the top ten PCA directions for the used dataset are shown in Fig. 3. As can be seen, the projection directions of a facial dataset are some facial features which are like ghost faces. That is why the facial projection directions are also referred to as *ghost faces*. The ghost faces in FDA and PCA are also referred to as *Fisherfaces* (Belhumeur et al., 1997; Etemad & Chellappa, 1997; Zhao et al., 1999) and *eigenfaces* (Turk &

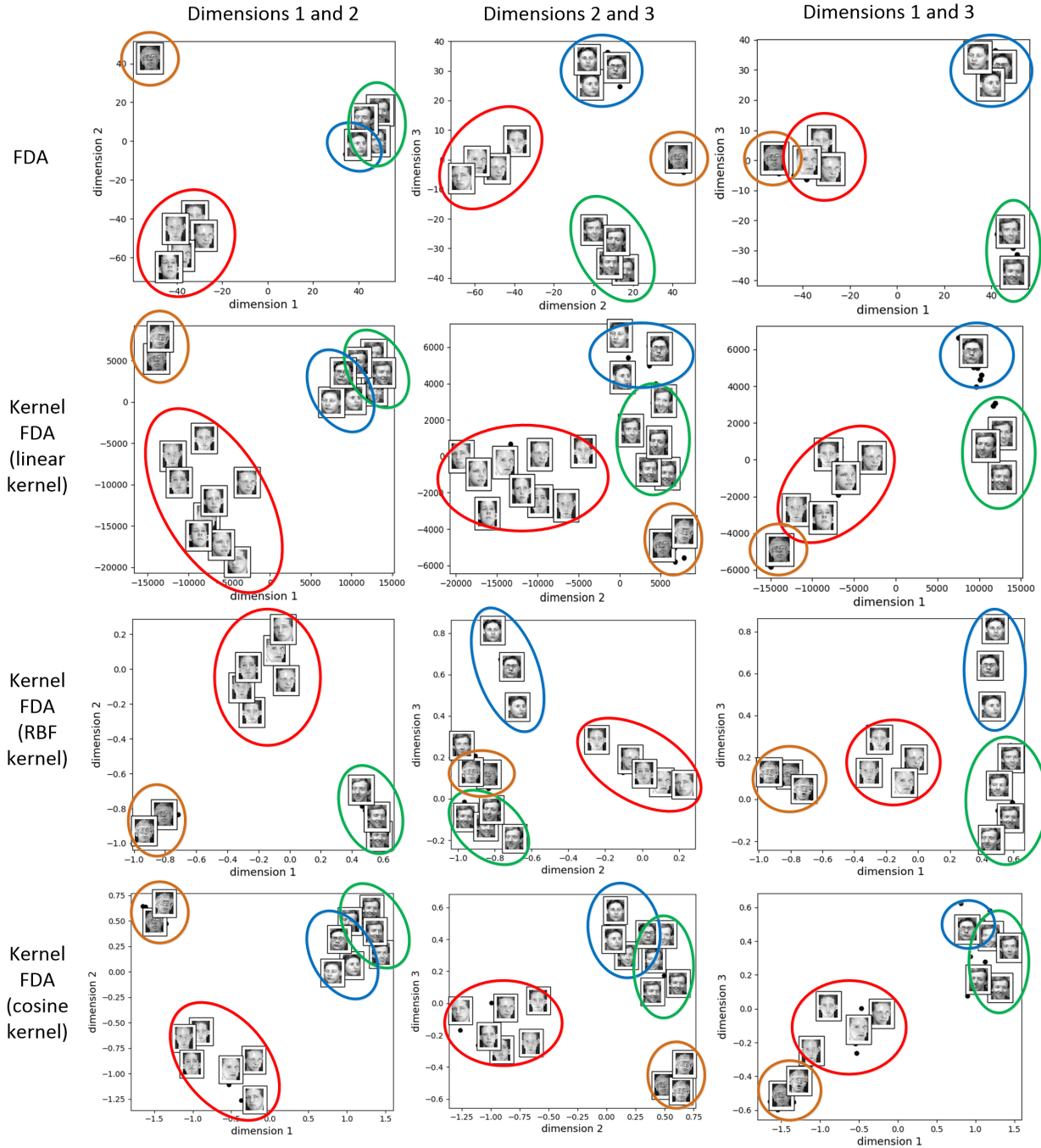


Figure 4. The projection of the first four classes of AT&T dataset onto FDA and kernel FDA subspaces where the used kernels were linear, RBF, and cosine kernels.

Pentland, 1991a;b), respectively. In Fig. 3, the projection directions have captured different facial features which discriminate the data with respect to the maximum variance in PCA and maximum class separation and minimum within class scatter in FDA. The captured features are eyes, eyeglasses, nose, cheeks, chin, lips, eyebrows, and hair, which are the most important facial features. This figure does

not include projection directions of kernel FDA and kernel PCA because in kernel FDA, the projection directions are n -dimensional and not d dimensional, and in kernel PCA, the projection directions are not available (see (Ghojogh & Crowley, 2019c)). Note that the face recognition using kernel FDA and kernel PCA are referred to as *kernel Fisherfaces* (Yang, 2002; Liu et al., 2004) and *kernel eigenfaces*

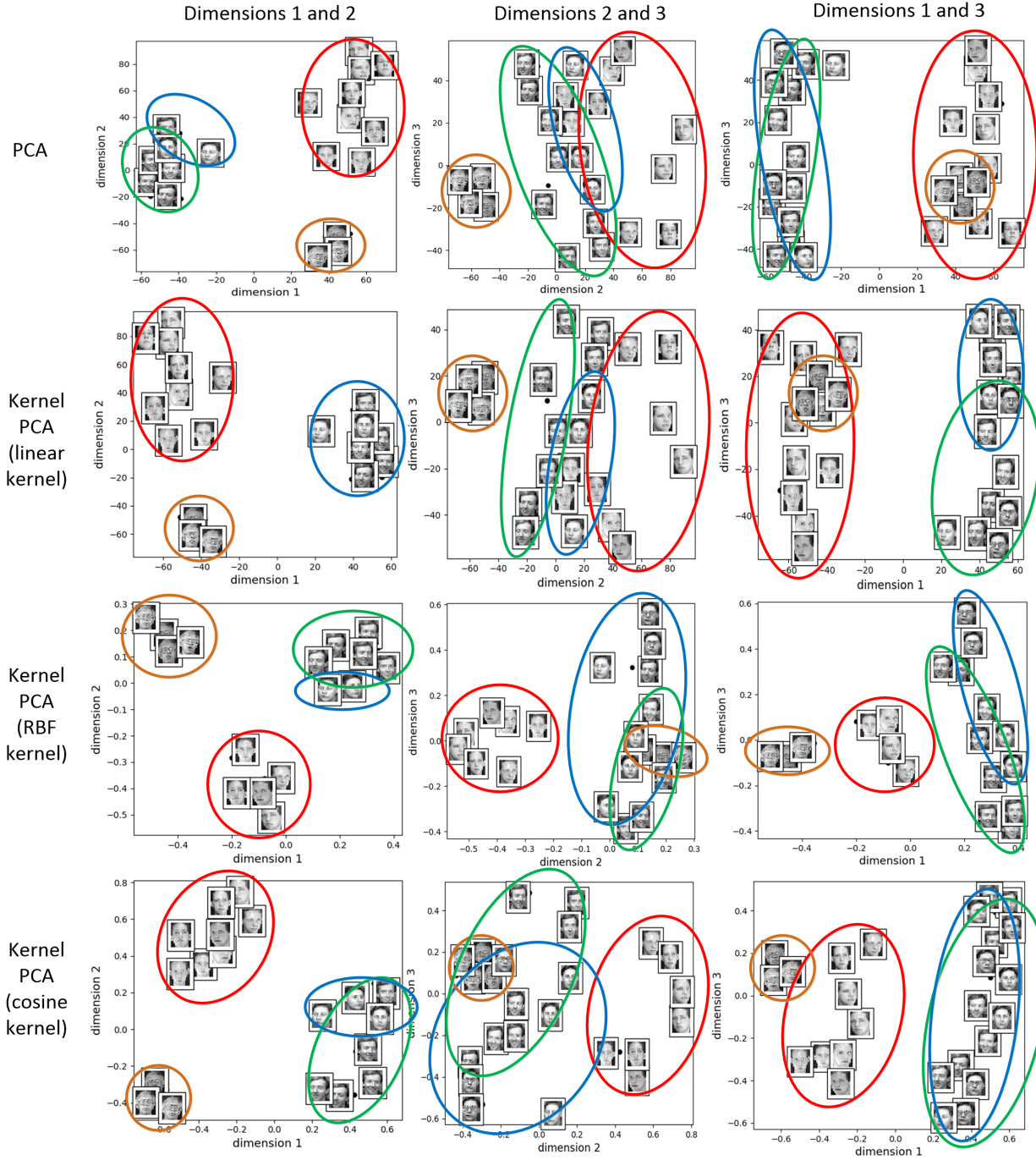


Figure 5. The projection of the first four classes of AT&T dataset onto PCA and kernel PCA subspaces where the used kernels were linear, RBF, and cosine kernels.

(Yang et al., 2000), respectively.

10.2. Projection of the Training Data

The projection of the images onto FDA and kernel FDA subspaces are shown in Fig. 4 where the linear, Radial Basis Function (RBF), and cosine kernels were used. The pro-

jection of the images using PCA and kernel PCA are also shown in Fig. 5. As can be seen, the FDA and kernel FDA subspaces have separated the classes much better than the PCA and kernel PCA subspaces. This is because the FDA and kernel FDA make use of the class labels in order to separate the classes in the subspace while the PCA and kernel

PCA only capture the variance (spread) of data regardless of class labels.

10.3. Reconstruction of Images

Figure 6 illustrates the reconstruction of some of the training images. The FDA has used its three projection directions for reconstruction. For reconstruction in PCA, we once used the top three PCA directions and one used the whole d PCA directions.

As can be seen in this figure, the reconstruction of PCA is much better than FDA. This makes sense because PCA is the best linear method for reconstruction having the least squared error (see (Ghojogh & Crowley, 2019c)). However, the responsibility of FDA is not reconstruction but separation of the classes. Thus, the FDA directions try to separate the classes as much as possible and do not *necessarily* care for a good reconstruction. Recall Fig. 2 which shows different cases for FDA and PCA directions. According to this figure, even in some datasets, FDA direction is orthogonal to PCA direction which is the best direction for reconstruction. It is noteworthy that reconstruction cannot be done in kernel FDA, so as in kernel PCA (see (Ghojogh & Crowley, 2019c)) as was mentioned before. Moreover, note that reconstruction can be done in FDA also for the out-of-sample data. Here, for the sake of brevity, we do not provide simulation for it.

10.4. Out-of-sample Projection

We took the first six images of each of the first four subjects in the AT&T dataset as the training images and the rest as the test (out-of-sample) images. The projection of the training and the out-of-sample images onto FDA and kernel FDA (using linear, RBF, and cosine kernels) are shown in Fig. 7. This figure shows that projection of out-of-sample images have been properly carried on in FDA and kernel FDA.

11. Conclusion

This paper was a tutorial paper introducing FDA and kernel FDA in detail. Various concepts about FDA, such as rank of scatters, dimensionality of the subspace, an example for interpretation, robust FDA, equivalency to LDA, and Fisher forest were explained and discussed. Both cases of two- and multi-classes were covered for FDA and kernel FDA. Finally, some simulations were performed to validate the theory in practice and compare to the unsupervised PCA method.

Acknowledgment

The authors hugely thank Prof. Ali Ghodsi (see his great online courses (Ghodsi, 2017; 2015)), Prof. Mu Zhu, Prof. Hoda Mohammadzade, and other professors whose courses have partly covered the materials mentioned in this tutorial

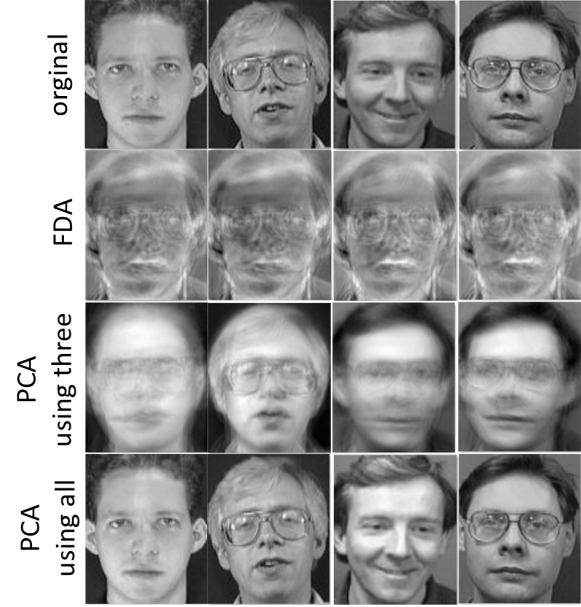


Figure 6. The reconstruction of four sample faces of AT&T datasets in FDA and PCA.

paper.

A. Metric Learning

The general form of metric (Peltonen et al., 2004) is usually defined as a form similar to Mahalanobis distance (McLachlan, 1999; De Maesschalck et al., 2000). The metric is:

$$\|x_i - x_j\|_A := (x_i - x_j)^T A (x_i - x_j), \quad (148)$$

where:

$$A = UU^T \succeq 0, \quad (149)$$

to have a valid distance metric. Most of the metric learning algorithms (Kulis et al., 2013) are optimization problems where A is unknown to make data points in same class (similar pairs) closer to each other, and points in different classes far apart from each other. We have:

$$\|x_i - x_j\|_A \stackrel{(149)}{=} (x_i - x_j)^T UU^T (x_i - x_j) \quad (150)$$

$$= (U^T x_i - U^T x_j)^T (U^T x_i - U^T x_j), \quad (151)$$

so this metric is equivalent to projection of data with projection matrix U and then using Euclidean distance in the embedded space (Peltonen et al., 2004). Therefore, Metric learning can be considered as a feature extraction (Ghojogh et al., 2019b) and manifold learning method (Alipanahi et al., 2008; Globerson & Roweis, 2006).

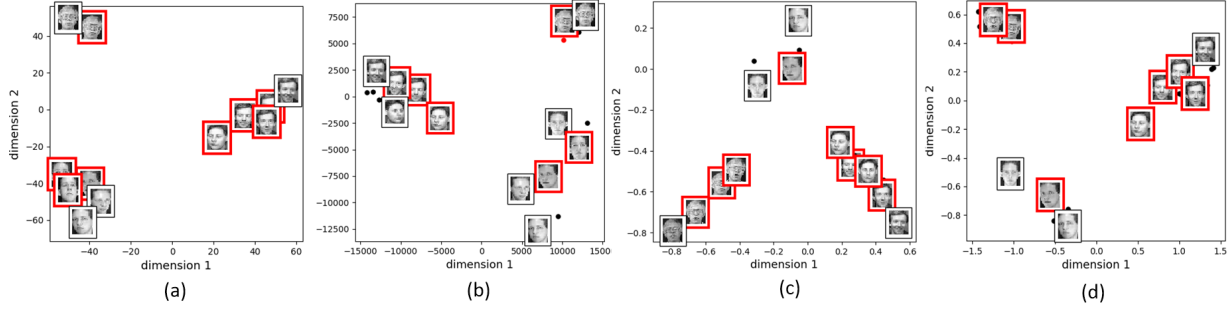


Figure 7. The first two dimensions of the projection of both training and out-of-sample instances in the first four classes of AT&T dataset onto subspaces of (a) FDA, (b) kernel FDA using linear kernel, (c) kernel FDA using RBF kernel, and (d) kernel FDA using cosine kernel.

B. Rayleigh-Ritz Quotient

The *Rayleigh-Ritz quotient* or *Rayleigh quotient* is defined as (Parlett, 1998; Croot, 2005):

$$\mathbb{R} \ni R(\mathbf{A}, \mathbf{x}) := \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}}, \quad (152)$$

where \mathbf{A} is a symmetric matrix and \mathbf{x} is a non-zero vector:

$$\mathbf{A} = \mathbf{A}^\top, \quad \mathbf{x} \neq \mathbf{0}. \quad (153)$$

One of the properties of the Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, c\mathbf{x}) = R(\mathbf{A}, \mathbf{x}), \quad (154)$$

where c is a scalar. The proof is that:

$$\begin{aligned} R(\mathbf{A}, c\mathbf{x}) &= \frac{(c\mathbf{x})^\top \mathbf{A} c\mathbf{x}}{(c\mathbf{x})^\top c\mathbf{x}} \stackrel{(a)}{=} \frac{c\mathbf{x}^\top \mathbf{A} c\mathbf{x}}{c\mathbf{x}^\top c\mathbf{x}} \\ &\stackrel{(b)}{=} \frac{c^2}{c^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} \stackrel{(152)}{=} R(\mathbf{A}, \mathbf{x}), \end{aligned}$$

where (a) and (b) are because c is a scalar.

Because of the Eq. (154), the optimization of the Rayleigh-Ritz quotient has an equivalent (Croot, 2005):

$$\begin{aligned} &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{x}) \stackrel{(a)}{=} \\ &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{x}) \stackrel{(b)}{=} \\ &\text{subject to} \quad \|\mathbf{x}\|_2 = 1, \\ &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad \mathbf{x}^\top \mathbf{A} \mathbf{x} \\ &\text{subject to} \quad \|\mathbf{x}\|_2 = 1, \end{aligned} \quad (155)$$

where (a) is because if we define $\mathbf{y} := (1/\|\mathbf{x}\|_2) \mathbf{x}$, the Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, \mathbf{y}) = \frac{\mathbf{y}^\top \mathbf{A} \mathbf{y}}{\mathbf{y}^\top \mathbf{y}} = \frac{1/\|\mathbf{x}\|_2^2}{1/\|\mathbf{x}\|_2^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{x}} = R(\mathbf{A}, \mathbf{x}), \quad (156)$$

and:

$$\|\mathbf{y}\|_2^2 = \frac{1}{\|\mathbf{x}\|_2^2} \times \|\mathbf{x}\|_2^2 = 1 \implies \|\mathbf{y}\|_2 = 1. \quad (157)$$

Thus, we have $R(\mathbf{A}, \mathbf{y})$ subject to $\|\mathbf{y}\|_2 = 1$. Changing the dummy variable \mathbf{y} to \mathbf{x} gives the Eq. (155). The (b) notices $\mathbf{x}^\top \mathbf{x} = 1$ because of the constraint $\|\mathbf{x}\|_2 = 1$.

Note that the constraint in Eq. (155) can be equal to any constant which is proved similarly. Moreover, note that the value of constant in the constraint is not important because it will be removed after taking derivative from the Lagrangian in optimization (Boyd & Vandenberghe, 2004). The *generalized Rayleigh-Ritz quotient* or *generalized Rayleigh quotient* is defined as (Parlett, 1998; Ghogh et al., 2019a):

$$\mathbb{R} \ni R(\mathbf{A}, \mathbf{B}; \mathbf{x}) := \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}}, \quad (158)$$

where \mathbf{A} and \mathbf{B} are symmetric matrices and \mathbf{x} is a non-zero vector:

$$\mathbf{A} = \mathbf{A}^\top, \quad \mathbf{B} = \mathbf{B}^\top, \quad \mathbf{x} \neq \mathbf{0}. \quad (159)$$

If the symmetric \mathbf{B} is positive definite:

$$\mathbf{B} \succ 0, \quad (160)$$

it has a Cholesky decomposition:

$$\mathbf{B} = \mathbf{C}\mathbf{C}^\top, \quad (161)$$

where \mathbf{C} is a lower triangular matrix. In case $\mathbf{B} \succ 0$, the generalized Rayleigh-Ritz quotient can be converted to a Rayleigh-Ritz quotient:

$$R(\mathbf{A}, \mathbf{B}; \mathbf{x}) = R(\mathbf{D}, \mathbf{C}^\top \mathbf{x}), \quad (162)$$

where:

$$\mathbf{D} := \mathbf{C}^{-1} \mathbf{A} \mathbf{C}^{-\top}. \quad (163)$$

The proof is:

$$\begin{aligned}
 \text{RHS} &= R(\mathbf{D}, \mathbf{C}^\top \mathbf{x}) \stackrel{(152)}{=} \frac{(\mathbf{C}^\top \mathbf{x})^\top \mathbf{D} (\mathbf{C}^\top \mathbf{x})}{(\mathbf{C}^\top \mathbf{x})^\top (\mathbf{C}^\top \mathbf{x})} \\
 &\stackrel{(163)}{=} \frac{\mathbf{x}^\top \mathbf{C} \mathbf{C}^{-1} \mathbf{A} (\mathbf{C} \mathbf{C}^{-1})^\top \mathbf{x}}{\mathbf{x}^\top (\mathbf{C} \mathbf{C}^\top) \mathbf{x}} \stackrel{(a)}{=} \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}} \\
 &\stackrel{(158)}{=} R(\mathbf{A}, \mathbf{B}; \mathbf{x}) = \text{LHS}, \quad \text{Q.E.D.},
 \end{aligned}$$

where RHS and LHS are short for right and left hand sides and (a) is because of Eq. (161) and $\mathbf{C} \mathbf{C}^{-1} = \mathbf{I}$ because \mathbf{C} is a square matrix.

Similarly, one of the properties of the generalized Rayleigh-Ritz quotient is:

$$R(\mathbf{A}, \mathbf{B}; c\mathbf{x}) = R(\mathbf{A}, \mathbf{B}; \mathbf{x}), \quad (164)$$

where c is a scalar. The proof is that:

$$\begin{aligned}
 R(\mathbf{A}, \mathbf{B}; c\mathbf{x}) &= \frac{(c\mathbf{x})^\top \mathbf{A} c\mathbf{x}}{(c\mathbf{x})^\top \mathbf{B} c\mathbf{x}} \stackrel{(a)}{=} \frac{c\mathbf{x}^\top \mathbf{A} c\mathbf{x}}{c\mathbf{x}^\top \mathbf{B} c\mathbf{x}} \\
 &\stackrel{(b)}{=} \frac{c^2}{c^2} \times \frac{\mathbf{x}^\top \mathbf{A} \mathbf{x}}{\mathbf{x}^\top \mathbf{B} \mathbf{x}} \stackrel{(158)}{=} R(\mathbf{A}, \mathbf{B}; \mathbf{x}),
 \end{aligned}$$

where (a) and (b) are because c is a scalar.

Because of the Eq. (164), the optimization of the generalized Rayleigh-Ritz quotient has an equivalent:

$$\begin{aligned}
 &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad R(\mathbf{A}, \mathbf{B}; \mathbf{x}) \equiv \\
 &\underset{\mathbf{x}}{\text{minimize/maximize}} \quad \mathbf{x}^\top \mathbf{A} \mathbf{x} \\
 &\text{subject to} \quad \mathbf{x}^\top \mathbf{B} \mathbf{x} = 1,
 \end{aligned} \quad (165)$$

for a similar reason that we provided for the Rayleigh-Ritz quotient. the constraint can be equal to any constant because in the derivative of Lagrangian, the constant will be dropped.

References

- Alipanahi, Babak, Biggs, Michael, and Ghodsi, Ali. Distance metric learning vs. Fisher discriminant analysis. In *Proceedings of the 23rd national conference on Artificial intelligence*, volume 2, pp. 598–603, 2008.
- Alperin, Jonathan L. *Local representation theory: Modular representations as an introduction to the local representation theory of finite groups*, volume 11. Cambridge University Press, 1993.
- Belhumeur, Peter N, Hespanha, João P, and Kriegman, David J. Eigenfaces vs. Fisherfaces: Recognition using class specific linear projection. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (7):711–720, 1997.
- Boyd, Stephen and Vandenberghe, Lieven. *Convex optimization*. Cambridge university press, 2004.
- Croot, Ernie. The Rayleigh principle for finding eigenvalues. Technical report, Georgia Institute of Technology, School of Mathematics, 2005. Online: http://people.math.gatech.edu/~ecroot/notes_linear.pdf, Accessed: March 2019.
- De Maesschalck, Roy, Jouan-Rimbaud, Delphine, and Massart, Désiré L. The Mahalanobis distance. *Chemometrics and intelligent laboratory systems*, 50(1):1–18, 2000.
- Deng, Weihong, Hu, Jiani, Guo, Jun, and Zhang, Hong-gang. Robust discriminant analysis of gabor feature for face recognition. In *Fourth International Conference on Fuzzy Systems and Knowledge Discovery (FSKD 2007)*, volume 3, pp. 248–252. IEEE, 2007.
- Etemad, Kamran and Chellappa, Rama. Discriminant analysis for recognition of human face images. *Journal of the Optical Society of America A*, 14(8):1724–1733, 1997.
- Fisher, Ronald A. Xv.the correlation between relatives on the supposition of mendelian inheritance. *Earth and Environmental Science Transactions of the Royal Society of Edinburgh*, 52(2):399–433, 1919.
- Fisher, Ronald A. The use of multiple measurements in taxonomic problems. *Annals of eugenics*, 7(2):179–188, 1936.
- Fisher, Ronald Aylmer. Statistical methods for research workers. In *Breakthroughs in statistics*, pp. 66–70. Springer, 1992.
- Frieden, B Roy. *Science from Fisher information: a unification*. Cambridge University Press, 2004.
- Friedman, Jerome, Hastie, Trevor, and Tibshirani, Robert. *The elements of statistical learning*, volume 2. Springer series in statistics New York, NY, USA, 2009.
- Ghodsi, Ali. Classification course, department of statistics and actuarial science, university of Waterloo. Online Youtube Videos, 2015. Accessed: January 2019.
- Ghodsi, Ali. Data visualization course, department of statistics and actuarial science, university of Waterloo. Online Youtube Videos, 2017. Accessed: January 2019.
- Ghojogh, Benyamin and Crowley, Mark. Linear and quadratic discriminant analysis: Tutorial. *arXiv preprint arXiv:1906.02590*, 2019a.
- Ghojogh, Benyamin and Crowley, Mark. The theory behind overfitting, cross validation, regularization, bagging, and boosting: Tutorial. *arXiv preprint arXiv:1905.12787*, 2019b.

- Ghojogh, Benyamin and Crowley, Mark. Unsupervised and supervised principal component analysis: Tutorial. *arXiv preprint arXiv:1906.03148*, 2019c.
- Ghojogh, Benyamin and Mohammadzade, Hoda. Automatic extraction of key-poses and key-joints for action recognition using 3d skeleton data. In *2017 10th Iranian Conference on Machine Vision and Image Processing (MVIP)*, pp. 164–170. IEEE, 2017.
- Ghojogh, Benyamin, Mohammadzade, Hoda, and Mokari, Mozghan. Fisherposes for human action recognition using kinect sensor data. *IEEE Sensors Journal*, 18(4): 1612–1627, 2017.
- Ghojogh, Benyamin, Karray, Fakhri, and Crowley, Mark. Eigenvalue and generalized eigenvalue problems: Tutorial. *arXiv preprint arXiv:1903.11240*, 2019a.
- Ghojogh, Benyamin, Samad, Maria N, Mashhadi, Sayema Asif, Kapoor, Tania, Ali, Wahab, Karray, Fakhri, and Crowley, Mark. Feature selection and feature extraction in pattern analysis: A literature review. *arXiv preprint arXiv:1905.02845*, 2019b.
- Globerson, Amir and Roweis, Sam T. Metric learning by collapsing classes. In *Advances in neural information processing systems*, pp. 451–458, 2006.
- Guo, Ming and Wang, Zhelong. A feature extraction method for human action recognition using body-worn inertial sensors. In *2015 IEEE 19th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, pp. 576–581. IEEE, 2015.
- Herbrich, Ralf. *Learning kernel classifiers: theory and algorithms*. Mit Press, 2001.
- Hofmann, Thomas, Schölkopf, Bernhard, and Smola, Alexander J. Kernel methods in machine learning. *The annals of statistics*, pp. 1171–1220, 2008.
- Kulis, Brian et al. Metric learning: A survey. *Foundations and Trends® in Machine Learning*, 5(4):287–364, 2013.
- Liu, Qingshan, Lu, Hanqing, and Ma, Songde. Improving kernel Fisher discriminant analysis for face recognition. *IEEE transactions on circuits and systems for video technology*, 14(1):42–49, 2004.
- McLachlan, Goeffrey J. Mahalanobis distance. *Resonance*, 4(6):20–26, 1999.
- Mika, Sebastian, Rätsch, Gunnar, Weston, Jason, Schölkopf, Bernhard, and Müller, Klaus-Robert. Fisher discriminant analysis with kernels. In *Proceedings of the 1999 IEEE signal processing society workshop on Neural networks for signal processing IX*, pp. 41–48. IEEE, 1999.
- Mika, Sebastian, Rätsch, Gunnar, Weston, Jason, Schölkopf, Bernhard, Smola, Alex J, and Müller, Klaus-Robert. Invariant feature extraction and classification in kernel spaces. In *Advances in neural information processing systems*, pp. 526–532, 2000.
- Mokari, Mozghan, Mohammadzade, Hoda, and Ghojogh, Benyamin. Recognizing involuntary actions from 3d skeleton data using body states. *Scientia Iranica*, 2018.
- Parlett, Beresford N. The symmetric eigenvalue problem. *Classics in Applied Mathematics*, 20, 1998.
- Peltonen, Jaakko, Klami, Arto, and Kaski, Samuel. Improved learning of Riemannian metrics for exploratory analysis. *Neural Networks*, 17(8-9):1087–1100, 2004.
- Polikar, Robi. Ensemble learning. In *Ensemble machine learning*, pp. 1–34. Springer, 2012.
- Samadani, Ali-Akbar, Ghodsi, Ali, and Kulić, Dana. Discriminative functional analysis of human movements. *Pattern Recognition Letters*, 34(15):1829–1839, 2013.
- Sugiyama, Masashi. Dimensionality reduction of multimodal labeled data by local fisher discriminant analysis. *Journal of machine learning research*, 8(May):1027–1061, 2007.
- Turk, Matthew and Pentland, Alex. Eigenfaces for recognition. *Journal of cognitive neuroscience*, 3(1):71–86, 1991a.
- Turk, Matthew A and Pentland, Alex P. Face recognition using eigenfaces. In *Computer Vision and Pattern Recognition, 1991. Proceedings CVPR'91., IEEE Computer Society Conference on*, pp. 586–591. IEEE, 1991b.
- Wang, Ruye. Generalized eigenvalue problem. <http://fourier.eng.hmc.edu/e161/lectures/algebra/node7.html>, 2015. Accessed: January 2019.
- Wang, Yanxia and Ruan, Qiuqi. Kernel fisher discriminant analysis for palmprint recognition. In *18th International Conference on Pattern Recognition (ICPR'06)*, volume 4, pp. 457–460. IEEE, 2006.
- Welling, Max. Fisher linear discriminant analysis. Technical report, Department of Computer Science, University of Toronto, 2005.
- Xu, Yong and Lu, Guangming. Analysis on Fisher discriminant criterion and linear separability of feature space. In *2006 International Conference on Computational Intelligence and Security*, volume 2, pp. 1671–1676. IEEE, 2006.

- Yang, M-H, Ahuja, Narendra, and Kriegman, David. Face recognition using kernel eigenfaces. In *Image processing, 2000. proceedings. 2000 international conference on*, volume 1, pp. 37–40. IEEE, 2000.
- Yang, Ming-Hsuan. Kernel Eigenfaces vs. kernel Fisher-faces: Face recognition using kernel methods. In *Proceedings of the fifth IEEE international conference on automatic face and gesture recognition*, pp. 215–220, 2002.
- Zhao, Wenyi, Chellappa, Rama, and Phillips, P Jonathon. *Subspace linear discriminant analysis for face recognition*. Citeseer, 1999.