
Local Fisher Discriminant Analysis for Supervised Dimensionality Reduction

Masashi Sugiyama

SUGI@CS.TITECH.AC.JP

Department of Computer Science, Tokyo Institute of Technology, 2-12-1-W8-74, O-okayama, Meguro-ku, Tokyo, 152-8552, Japan

Abstract

Dimensionality reduction is one of the important preprocessing steps in high-dimensional data analysis. In this paper, we consider the supervised dimensionality reduction problem where samples are accompanied with class labels. Traditional Fisher discriminant analysis is a popular and powerful method for this purpose. However, it tends to give undesired results if samples in some class form several separate clusters, i.e., multimodal. In this paper, we propose a new dimensionality reduction method called local Fisher discriminant analysis (LFDA), which is a localized variant of Fisher discriminant analysis. LFDA takes local structure of the data into account so the multimodal data can be embedded appropriately. We also show that LFDA can be extended to non-linear dimensionality reduction scenarios by the kernel trick.

1. Introduction

The goal of dimensionality reduction is to embed high-dimensional data samples in a low-dimensional space while most of ‘intrinsic information’ contained in the data is preserved. Once dimensionality reduction is carried out appropriately, we can utilize the compact representation of the data for various succeeding tasks such as visualization, classification, etc. In this paper, we consider the supervised dimensionality reduction problem where samples are accompanied with class labels.

Fisher discriminant analysis (FDA) (Fisher, 1936; Fukunaga, 1990) is a popular method for linear dimensionality reduction¹, which maximizes between-class scatter and minimizes within-class scatter. This traditional FDA is known to work well and is practically useful even now. However, it tends to give undesired results if samples in some class form several separate clusters (i.e., *multimodal*) (Fukunaga, 1990).

Multimodality is often observed in many practical applications. For example, in disease diagnosis, the distribution of medical checkup samples of sick patients could be multimodal since there may be several different causes even for a particular disease. Even in a traditional task of hand-written digit recognition, multimodality appears if digits are classified into, e.g., even and odd numbers. More generally, multimodality can naturally appear if multi-class classification problems are solved by a set of two-class ‘one-versus-rest’ problems.

To embed multimodal data well, it is important to preserve local structure of the data. *Locality-preserving projection* (LPP) (He & Niyogi, 2004) meets this requirement. LPP keeps nearby data pairs in the original space close in the embedding space, by which multimodal data can be embedded without losing its local structure. However, LPP is an unsupervised dimensionality reduction method which does not take the label information into account. Therefore, it does not necessarily work appropriately in supervised dimensionality reduction scenarios.

In this paper, we propose a new dimensionality reduction method called *local Fisher discriminant analysis*

¹Usually, FDA may refer to a classification method which first projects the data samples onto a one-dimensional space and then classifies the samples by thresholding. The one-dimensional embedding space used above is given as the maximizer of the so-called *Fisher criterion*. This Fisher criterion is often used for dimensionality reduction to a subspace with dimension more than one (Fukunaga, 1990). With some abuse, we refer to the dimensionality reduction method based on the Fisher criterion as FDA in the following (see Section 2.2 for detail).

(LFDA). LFDA combines the ideas of FDA and LPP: between-class separability is maximized while *within-class local structure* is preserved.

Because of the local structure preservation property, multimodal labeled data can be effectively embedded by LFDA. Furthermore, LFDA can give more separate embedding than FDA. To explain the reason intuitively, we first note that FDA can be regarded as maximizing between-class scatter under the constraint of keeping within-class scatter to a certain level (Fukunaga, 1990). When samples in some class are multimodal, keeping within-class scatter to a certain level appears to be quite hard since multimodal samples should be typically merged into a single cluster. Due to this strong constraint, less degree of freedom is left for increasing separability, and thus FDA results in less separate embedding. On the other hand, LFDA does not require multimodal samples to fall into a single cluster. As a result, more degree of freedom is left for increasing separability and thus highly separate embedding can be obtained.

By the so-called *kernel trick* (Schölkopf & Smola, 2002), FDA and LPP can be extended to non-linear dimensionality reduction scenarios (Mika et al., 2003; Belkin & Niyogi, 2003). We show that LFDA can also be non-linearized by the kernel trick.

2. Linear Dimensionality Reduction

In this section, we formulate the problem of linear dimensionality reduction and review typical existing methods.

2.1. Formulation

Let $\mathbf{x}_i \in \mathbb{R}^d$ ($i = 1, 2, \dots, n$) be d -dimensional samples and $y_i \in \{1, 2, \dots, \ell\}$ be associated class labels, where n is the number of samples and ℓ is the number of classes. Let n_i be the number of samples in the class i :

$$\sum_{i=1}^{\ell} n_i = n. \quad (1)$$

Let \mathbf{X} be the matrix of all samples:

$$\mathbf{X} = (\mathbf{x}_1 | \mathbf{x}_2 | \dots | \mathbf{x}_n). \quad (2)$$

Let $\mathbf{z}_i \in \mathbb{R}^m$ ($1 \leq m \leq d$) be embedded samples, where m is the dimension of the embedding space. Effectively we consider d to be large and m to be small, but not limited to such cases.

For the moment, we focus on linear dimensionality reduction, i.e., using a $d \times m$ transformation matrix \mathbf{T} , \mathbf{z}_i is given by

$$\mathbf{z}_i = \mathbf{T}^\top \mathbf{x}_i. \quad (3)$$

Later in Section 3.5, we discuss non-linear dimensionality reduction scenarios.

2.2. Fisher Linear Discriminant

Here we briefly review the definition of Fisher criterion for dimensionality reduction (Fisher, 1936; Fukunaga, 1990). With some abuse, we refer to the dimensionality reduction method based on the Fisher criterion as Fisher discriminant analysis (FDA), where the original FDA embeds the data samples only in one-dimensional space.

Let $\mathbf{S}^{(w)}$ and $\mathbf{S}^{(b)}$ be the *within-class scatter matrix* and the *between-class scatter matrix* defined by

$$\mathbf{S}^{(w)} = \sum_{i=1}^{\ell} \sum_{j: y_j=i} (\mathbf{x}_j - \boldsymbol{\mu}_i)(\mathbf{x}_j - \boldsymbol{\mu}_i)^\top, \quad (4)$$

$$\mathbf{S}^{(b)} = \sum_{i=1}^{\ell} n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^\top, \quad (5)$$

where $^\top$ denotes the transpose, $\boldsymbol{\mu}_i$ is the mean of samples in the class i and $\boldsymbol{\mu}$ is the mean of all samples:

$$\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{j: y_j=i} \mathbf{x}_j, \quad (6)$$

$$\boldsymbol{\mu} = \frac{1}{n} \sum_{i=1}^n \mathbf{x}_i. \quad (7)$$

Using $\mathbf{S}^{(w)}$ and $\mathbf{S}^{(b)}$, the FDA transformation matrix \mathbf{T}_{FDA} is defined as follows:

$$\mathbf{T}_{FDA} = \operatorname{argmax}_{\mathbf{T} \in \mathbb{R}^{d \times m}} \operatorname{tr} \left((\mathbf{T}^\top \mathbf{S}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \mathbf{S}^{(b)} \mathbf{T} \right). \quad (8)$$

That is, \mathbf{T} is determined so that between-class scatter is maximized while within-class scatter is minimized.

It is known that \mathbf{T}_{FDA} is given by

$$\mathbf{T}_{FDA} = (\boldsymbol{\varphi}_1 | \boldsymbol{\varphi}_2 | \dots | \boldsymbol{\varphi}_m), \quad (9)$$

where $\{\boldsymbol{\varphi}_i\}_{i=1}^d$ are the generalized eigenvectors associated to the generalized eigenvalues $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d$ of the following generalized eigenvalue problem:

$$\mathbf{S}^{(b)} \boldsymbol{\varphi} = \lambda \mathbf{S}^{(w)} \boldsymbol{\varphi}. \quad (10)$$

2.3. Locality-Preserving Projection

Here we briefly review the definition of locality-preserving projection (LPP) (He & Niyogi, 2004). Let \mathbf{A} be the *affinity matrix*, i.e., the n -dimensional matrix with the (i, j) -th element being the affinity between \mathbf{x}_i and \mathbf{x}_j . We suppose that elements of \mathbf{A} are in $[0, 1]$ and take larger values if \mathbf{x}_i and \mathbf{x}_j are ‘close’.

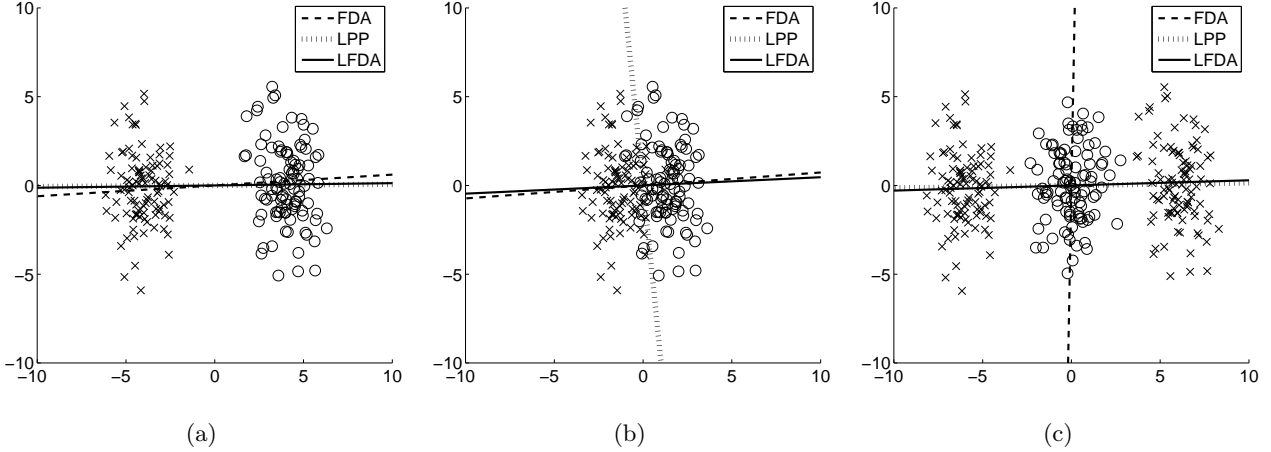


Figure 1. Examples of dimensionality reduction by FDA, LPP and LFDA.

There are several different manners in defining \mathbf{A} . A simple one is to define $\mathbf{A}_{i,j} = 1$ if \mathbf{x}_j is the k -nearest neighbor of \mathbf{x}_i or vice versa; otherwise $\mathbf{A}_{i,j} = 0$. More elaborate methods can be found, e.g., in (He & Niyogi, 2004; Zelnik-Manor & Perona, 2005). The latter may be useful for those who do not have any subjective/prior preference.

Using \mathbf{A} , the LPP transformation matrix \mathbf{T}_{LPP} is defined as follows.

$$\mathbf{T}_{LPP} = \underset{\mathbf{T} \in \mathbb{R}^{d \times m}}{\operatorname{argmin}} \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j} \|\mathbf{T}^\top \mathbf{x}_i - \mathbf{T}^\top \mathbf{x}_j\|^2$$

subject to $\mathbf{T}^\top \mathbf{X} \mathbf{D} \mathbf{X}^\top \mathbf{T} = \mathbf{I}$, (11)

where \mathbf{I} is the identity matrix and \mathbf{D} is the n -dimensional diagonal matrix with i -th diagonal element being

$$D_{i,i} = \sum_{j=1}^n \mathbf{A}_{i,j}. \quad (12)$$

Eq.(11) implies that in LPP, \mathbf{T} is determined so that *nearby* data pairs in the original space are kept close in the embedding space. Note that $\mathbf{T}^\top \mathbf{X} \mathbf{D} \mathbf{X}^\top \mathbf{T} = \mathbf{I}$ is a constraint to avoid a trivial solution $\mathbf{T} = \mathbf{O}$.

It is known that the LPP transformation matrix \mathbf{T}_{LPP} is given by

$$\mathbf{T}_{LPP} = (\psi_{d-m+1} | \psi_{d-m+2} | \cdots | \psi_d), \quad (13)$$

where $\{\psi_i\}_{i=1}^d$ are the eigenvectors associated to the eigenvalues $\gamma_1 \geq \gamma_2 \geq \cdots \geq \gamma_d$ of the following eigenvalue problem:

$$\mathbf{X} \mathbf{L} \mathbf{X}^\top \psi = \gamma \mathbf{X} \mathbf{D} \mathbf{X}^\top \psi, \quad (14)$$

where

$$\mathbf{L} = \mathbf{D} - \mathbf{A}. \quad (15)$$

2.4. Typical Behavior

In Figure 1, dimensionality reduction results obtained by FDA and LPP are shown, where two-dimensional two-class data samples are embedded in a one-dimensional subspace. In LPP, the affinity matrix \mathbf{A} is determined by the *local scaling method* (Zelnik-Manor & Perona, 2005).

For the simplest data set depicted in Figure 1(a), both FDA and LPP give reasonable results where samples of different classes are nicely separated. For the data set depicted in Figure 1(b), FDA still works well. However, LPP mixes samples of different classes into one cluster, which is caused by the unsupervised nature of LPP. On the other hand, for the data set depicted in Figure 1(c), LPP works well but FDA gives an undesired result. The reason is that the levels of between-class scatter and within-class scatter are not evaluated in an intuitively natural way because of the two separate clusters of the same class. See also (Fukunaga, 1990).

3. Local Fisher Discriminant Analysis

As illustrated above, FDA can perform poorly if samples in some class form several separate clusters (i.e., multimodal). In other words, the undesired behavior is caused by the *globality* when evaluating within-class scatter and between-class scatter. On the other hand, LPP can make samples of different classes overlapped if they are close in the original high-dimensional space.

To overcome this problem, we propose combining the idea of FDA and LPP: we evaluate the levels of within-class scatter and between-class scatter in a *local* manner, by which class separability and local structure

preservation could be attained at the same time. We call our new method *local Fisher discriminant analysis* (LFDA).

3.1. Reformulating FDA

To introduce our new method, let us first reformulate FDA in a *pairwise* manner.

Lemma 1 $\mathbf{S}^{(w)}$ and $\mathbf{S}^{(b)}$ are expressed as

$$\mathbf{S}^{(w)} = \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(w)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \quad (16)$$

$$\mathbf{S}^{(b)} = \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(b)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \quad (17)$$

where

$$\mathbf{A}_{i,j}^{(w)} = \begin{cases} 1/n_c & \text{if } y_i = y_j = c, \\ 0 & \text{if } y_i \neq y_j, \end{cases} \quad (18)$$

$$\mathbf{A}_{i,j}^{(b)} = \begin{cases} 1/n - 1/n_c & \text{if } y_i = y_j = c, \\ 1/n & \text{if } y_i \neq y_j, \end{cases} \quad (19)$$

(Proof) It follows from Eq.(4) that

$$\begin{aligned} \mathbf{S}^{(w)} &= \sum_{i=1}^{\ell} \sum_{j:y_j=i} \left(\mathbf{x}_j - \frac{1}{n_i} \sum_{p:y_p=i} \mathbf{x}_p \right) \left(\mathbf{x}_j - \frac{1}{n_i} \sum_{q:y_q=i} \mathbf{x}_q \right)^\top \\ &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i=1}^{\ell} \frac{1}{n_i} \sum_{p,q:y_p=y_q=i} \mathbf{x}_p \mathbf{x}_q^\top \\ &= \sum_{i=1}^n \left(\sum_{j=1}^n \mathbf{A}_{i,j}^{(w)} \right) \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(w)} \mathbf{x}_i \mathbf{x}_j^\top \\ &= \frac{1}{2} \sum_{i,j=1}^n \mathbf{A}_{i,j}^{(w)} (\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top), \end{aligned} \quad (20)$$

which yields Eq.(16). Let $\mathbf{S}^{(m)}$ be the *mixture scatter matrix* (Fukunaga, 1990):

$$\begin{aligned} \mathbf{S}^{(m)} &= \mathbf{S}^{(w)} + \mathbf{S}^{(b)} \\ &= \sum_{i=1}^n (\mathbf{x}_i - \boldsymbol{\mu})(\mathbf{x}_i - \boldsymbol{\mu})^\top. \end{aligned} \quad (21)$$

Then we have

$$\begin{aligned} \mathbf{S}^{(b)} &= \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^\top - \frac{1}{n} \sum_{i,j=1}^n \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{S}^{(w)} \\ &= \sum_{i=1}^n \left(\sum_{j=1}^n \frac{1}{n} \right) \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i,j=1}^n \frac{1}{n} \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{S}^{(w)} \\ &= \frac{1}{2} \sum_{i,j=1}^n \left(\frac{1}{n} - \mathbf{A}_{i,j}^{(w)} \right) \\ &\quad \times (\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top), \end{aligned} \quad (22)$$

which yields Eq.(17). \blacksquare

Note that $1/n - 1/n_c$ is negative while $1/n_c$ and $1/n$ are positive.

The above pairwise representation gives us a new interpretation of FDA: it tries to keep in-class data pairs close (since $\mathbf{A}_{i,j}^{(w)}$ is positive and $\mathbf{A}_{i,j}^{(b)}$ is negative if $y_i = y_j$) and between-class data pairs apart (since $\mathbf{A}_{i,j}^{(b)}$ is positive if $y_i \neq y_j$).

3.2. Definition of LFDA

Based on the above pairwise expression, let us define the *local* within-class scatter matrix $\bar{\mathbf{S}}^{(w)}$ and the *local* between-class scatter matrix $\bar{\mathbf{S}}^{(b)}$ as follows.

$$\bar{\mathbf{S}}^{(w)} = \frac{1}{2} \sum_{i,j=1}^n \bar{\mathbf{A}}_{i,j}^{(w)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \quad (23)$$

$$\bar{\mathbf{S}}^{(b)} = \frac{1}{2} \sum_{i,j=1}^n \bar{\mathbf{A}}_{i,j}^{(b)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \quad (24)$$

where

$$\bar{\mathbf{A}}_{i,j}^{(w)} = \begin{cases} \mathbf{A}_{i,j}/n_c & \text{if } y_i = y_j = c, \\ 0 & \text{if } y_i \neq y_j, \end{cases} \quad (25)$$

$$\bar{\mathbf{A}}_{i,j}^{(b)} = \begin{cases} \mathbf{A}_{i,j}(1/n - 1/n_c) & \text{if } y_i = y_j = c, \\ 1/n & \text{if } y_i \neq y_j. \end{cases} \quad (26)$$

Compared with the global counterparts $\mathbf{S}^{(w)}$ and $\mathbf{S}^{(b)}$, the values for in-class pairs are weighted by the affinity in $\bar{\mathbf{S}}^{(w)}$ and $\bar{\mathbf{S}}^{(b)}$. This means that far-apart in-class pairs have less influence in $\bar{\mathbf{S}}^{(w)}$ and $\bar{\mathbf{S}}^{(b)}$.

Using $\bar{\mathbf{S}}^{(w)}$ and $\bar{\mathbf{S}}^{(b)}$, we define the LFDA transformation matrix \mathbf{T}_{LFDA} as

$$\mathbf{T}_{LFDA} = \operatorname{argmax}_{\mathbf{T} \in \mathbb{R}^{d \times m}} \operatorname{tr} \left((\mathbf{T}^\top \bar{\mathbf{S}}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \bar{\mathbf{S}}^{(b)} \mathbf{T} \right). \quad (27)$$

That is, we determine \mathbf{T} so that *nearby* data pairs of the *same class* are close and data pairs of different

classes are apart. Data pairs of the same class but far apart are not imposed to be close.

If the affinity matrix \mathbf{A} is taken to be one for all in-class pairs (i.e., all in-class pairs are ‘equally close’ each other), $\bar{\mathbf{S}}^{(w)}$ and $\bar{\mathbf{S}}^{(b)}$ are reduced to $\mathbf{S}^{(w)}$ and $\mathbf{S}^{(b)}$ so LFDA is reduced to the ordinary FDA. Therefore, LFDA may be regarded as a natural localized variant of FDA.

Since Eq.(27) is the same form as Eq.(8), we can easily obtain \mathbf{T}_{LFDA} by solving a generalized eigenvalue problem. In Section 3.4, we show an efficient method for computing \mathbf{T}_{LFDA} .

Examples of dimensionality reduction by LFDA are illustrated in Figure 1, where the affinity matrix \mathbf{A} is determined by the same way as being done for LPP. The figures show that LFDA gives desirable results for all three data sets, i.e., LFDA can compensate the drawbacks of FDA and LPP.

3.3. Why is LFDA good?

LFDA does not impose far-apart data pairs of the same class to be close, by which local structure of the data tends to be preserved. This itself is a useful property, e.g., in data visualization tasks because ‘interesting’ structure such as multimodality is not lost by dimensionality reduction.

In addition to that, LFDA can even provide more separate embedding than FDA. To explain the reason, let us recall that the FDA transformation matrix \mathbf{T}_{FDA} can also be obtained as follows (Fukunaga, 1990).

$$\begin{aligned} \mathbf{T}_{FDA} = \operatorname{argmax}_{\mathbf{T} \in \mathbb{R}^{d \times m}} \operatorname{tr} \left(\mathbf{T}^\top \mathbf{S}^{(b)} \mathbf{T} \right) \\ \text{subject to } \mathbf{T}^\top \mathbf{S}^{(w)} \mathbf{T} = \mathbf{I}. \end{aligned} \quad (28)$$

This representation implies that FDA maximizes between-class scatter under the constraint of keeping within-class scatter to a certain level.

When one of the classes is multimodal, the above constraint is actually quite restrictive since the multimodal data should be typically merged into a single cluster. Therefore, the remaining degree of freedom which can be used for maximizing between-class scatter may not be that much. As a result, FDA can result in less separate embedding.

On the other hand, it is straightforward to show that \mathbf{T}_{LFDA} can also be expressed as follows.

$$\begin{aligned} \mathbf{T}_{LFDA} = \operatorname{argmax}_{\mathbf{T} \in \mathbb{R}^{d \times m}} \operatorname{tr} \left(\mathbf{T}^\top \bar{\mathbf{S}}^{(b)} \mathbf{T} \right) \\ \text{subject to } \mathbf{T}^\top \bar{\mathbf{S}}^{(w)} \mathbf{T} = \mathbf{I}. \end{aligned} \quad (29)$$

The constraint in Eq.(29) is less restrictive than that in Eq.(28) since faraway pairs of in-class samples are not imposed to be close. Due to this weaker constraint, the degree of freedom which can be used for maximizing separability between different classes remains more than FDA. Thus, LFDA can result in more separate embedding.

3.4. Efficiently Computing LFDA Transformation Matrix

Here, we provide an efficient method to compute the LFDA transformation matrix \mathbf{T}_{LFDA} .

Let $\bar{\mathbf{S}}^{(m)}$ be the *local* mixture scatter matrix defined by

$$\begin{aligned} \bar{\mathbf{S}}^{(m)} &= \bar{\mathbf{S}}^{(w)} + \bar{\mathbf{S}}^{(b)} \\ &= \frac{1}{2} \sum_{i,j=1}^n \bar{\mathbf{A}}_{i,j}^{(m)} (\mathbf{x}_i - \mathbf{x}_j)(\mathbf{x}_i - \mathbf{x}_j)^\top, \end{aligned} \quad (30)$$

where $\bar{\mathbf{A}}^{(m)}$ is the n -dimensional matrix with (i, j) -th element being

$$\bar{\mathbf{A}}_{i,j}^{(m)} = \bar{\mathbf{A}}_{i,j}^{(w)} + \bar{\mathbf{A}}_{i,j}^{(b)} = \begin{cases} \mathbf{A}_{i,j}/n & \text{if } y_i = y_j, \\ 1/n & \text{if } y_i \neq y_j. \end{cases} \quad (31)$$

Then the same solution \mathbf{T}_{LFDA} can be still obtained if $\bar{\mathbf{S}}^{(b)}$ in Eq.(27) is replaced by $\bar{\mathbf{S}}^{(m)}$ because of the following identity (cf. (Fukunaga, 1990)):

$$\begin{aligned} \operatorname{tr} \left((\mathbf{T}^\top \bar{\mathbf{S}}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \bar{\mathbf{S}}^{(m)} \mathbf{T} \right) \\ = \operatorname{tr} \left((\mathbf{T}^\top \bar{\mathbf{S}}^{(w)} \mathbf{T})^{-1} \mathbf{T}^\top \bar{\mathbf{S}}^{(b)} \mathbf{T} \right) + m. \end{aligned} \quad (32)$$

In the following, we employ the criterion with $\bar{\mathbf{S}}^{(m)}$ since it is simpler.

The left-hand side of Eq.(32) is the same form as FDA. Therefore, the solution \mathbf{T}_{LFDA} is analytically given as follows.

$$\mathbf{T}_{LFDA} = (\bar{\varphi}_1 | \bar{\varphi}_2 | \cdots | \bar{\varphi}_m), \quad (33)$$

where $\{\bar{\varphi}_i\}_{i=1}^d$ are the generalized eigenvectors associated to the generalized eigenvalues $\bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \cdots \geq \bar{\lambda}_d$ of the following generalized eigenvalue problem:

$$\bar{\mathbf{S}}^{(m)} \bar{\varphi} = \bar{\lambda} \bar{\mathbf{S}}^{(w)} \bar{\varphi}. \quad (34)$$

Below, we show how $\bar{\mathbf{S}}^{(m)}$ and $\bar{\mathbf{S}}^{(w)}$ can be computed efficiently.

Since

$$\begin{aligned}\bar{\mathbf{S}}^{(m)} &= \frac{1}{2} \sum_{i,j=1}^n \bar{\mathbf{A}}_{i,j}^{(m)} (\mathbf{x}_i \mathbf{x}_i^\top + \mathbf{x}_j \mathbf{x}_j^\top - \mathbf{x}_i \mathbf{x}_j^\top - \mathbf{x}_j \mathbf{x}_i^\top) \\ &= \sum_{i=1}^n \left(\sum_{j=1}^n \bar{\mathbf{A}}_{i,j}^{(m)} \right) \mathbf{x}_i \mathbf{x}_i^\top - \sum_{i,j=1}^n \bar{\mathbf{A}}_{i,j}^{(m)} \mathbf{x}_i \mathbf{x}_j^\top,\end{aligned}\quad (35)$$

$\bar{\mathbf{S}}^{(m)}$ can be expressed in a matrix form as

$$\bar{\mathbf{S}}^{(m)} = \mathbf{X} \bar{\mathbf{L}}^{(m)} \mathbf{X}^\top, \quad (36)$$

where $\bar{\mathbf{L}}^{(m)} = \bar{\mathbf{D}}^{(m)} - \bar{\mathbf{A}}^{(m)}$ and $\bar{\mathbf{D}}^{(m)}$ is the n -dimensional diagonal matrix with i -th diagonal element being $\bar{D}_{i,i}^{(m)} = \sum_{j=1}^n \bar{\mathbf{A}}_{i,j}^{(m)}$. Similarly, $\bar{\mathbf{S}}^{(w)}$ can be expressed in a matrix form as

$$\bar{\mathbf{S}}^{(w)} = \mathbf{X} \bar{\mathbf{L}}^{(w)} \mathbf{X}^\top, \quad (37)$$

where $\bar{\mathbf{L}}^{(w)} = \bar{\mathbf{D}}^{(w)} - \bar{\mathbf{A}}^{(w)}$ and $\bar{\mathbf{D}}^{(w)}$ is the n -dimensional diagonal matrix with i -th diagonal element being $\bar{D}_{i,i}^{(w)} = \sum_{j=1}^n \bar{\mathbf{A}}_{i,j}^{(w)}$.

$\bar{\mathbf{L}}^{(m)}$ and $\bar{\mathbf{L}}^{(w)}$ are n -dimensional matrices so could be very high dimensional. However, $\bar{\mathbf{L}}^{(w)}$ could be *sparse* so computing $\bar{\mathbf{S}}^{(w)}$ directly by Eq.(37) may be efficient. On the other hand, computing $\bar{\mathbf{S}}^{(m)}$ directly by Eq.(36) is not so efficient since $\bar{\mathbf{A}}^{(m)}$ is *dense*.

This problem can be alleviated as follows. $\bar{\mathbf{A}}^{(m)}$ can be decomposed as

$$\bar{\mathbf{A}}^{(m)} = \mathbf{1}\mathbf{1}^\top/n + \bar{\mathbf{A}}^{(m')}, \quad (38)$$

where $\mathbf{1}$ is the n -dimensional vector with all ones and $\bar{\mathbf{A}}^{(m')}$ is the n -dimensional matrix with (i,j) -th element being

$$\bar{A}_{i,j}^{(m')} = \begin{cases} (\mathbf{A}_{i,j} - 1)/n & \text{if } y_i = y_j, \\ 0 & \text{if } y_i \neq y_j. \end{cases} \quad (39)$$

Then $\bar{\mathbf{S}}^{(m)}$ can be expressed as follows:

$$\bar{\mathbf{S}}^{(m)} = \mathbf{X} \bar{\mathbf{D}}^{(m)} \mathbf{X}^\top - (\mathbf{X}\mathbf{1})(\mathbf{X}\mathbf{1})^\top/n - \mathbf{X} \bar{\mathbf{A}}^{(m')} \mathbf{X}^\top, \quad (40)$$

where $\bar{\mathbf{D}}^{(m)}$ is expressed by using $\bar{\mathbf{A}}^{(m')}$ as

$$\bar{D}_{i,i}^{(m)} = 1 + \sum_{j=1}^n \bar{\mathbf{A}}_{i,j}^{(m')}. \quad (41)$$

$\bar{\mathbf{A}}^{(m')}$ becomes a block-diagonal matrix if $\{\mathbf{x}_i\}_{i=1}^n$ are sorted by the labels, which implies that the third term

in the right-hand side of Eq.(40) may be computed efficiently. Since the first two terms in the right-hand side of Eq.(40) can also be computed efficiently, computing $\bar{\mathbf{S}}^{(m)}$ by Eq.(40) may be more efficient than directly by Eq.(36).

To further improve computational efficiency, the affinity matrix \mathbf{A} may be computed in a classwise manner, i.e., the between-class affinity is set to be zero.

3.5. Kernel LFDA for Non-Linear Dimensionality Reduction

So far, we focused on linear dimensionality reduction. Here we extend our discussion to non-linear dimensionality reduction scenarios.

From Eqs.(36) and (37), the generalized eigenvalue problem (34) can be expressed as

$$\mathbf{X} \bar{\mathbf{L}}^{(m)} \mathbf{X}^\top \bar{\boldsymbol{\varphi}} = \bar{\lambda} \mathbf{X} \bar{\mathbf{L}}^{(w)} \mathbf{X}^\top \bar{\boldsymbol{\varphi}}. \quad (42)$$

When $d \leq n$, any vector $\bar{\boldsymbol{\varphi}} \in \mathbb{R}^d$ can be expressed by using some vector $\bar{\boldsymbol{\alpha}} \in \mathbb{R}^n$ as $\bar{\boldsymbol{\varphi}} = \mathbf{X} \bar{\boldsymbol{\alpha}}$. Then, multiplying Eq.(42) by \mathbf{X}^\top from the left-hand side, we have

$$\mathbf{K} \bar{\mathbf{L}}^{(m)} \mathbf{K} \bar{\boldsymbol{\alpha}} = \bar{\lambda} \mathbf{K} \bar{\mathbf{L}}^{(w)} \mathbf{K} \bar{\boldsymbol{\alpha}}, \quad (43)$$

where \mathbf{K} is the n -dimensional matrix with the (i,j) -th element being

$$\mathbf{K}_{i,j} = \mathbf{x}_i^\top \mathbf{x}_j. \quad (44)$$

This implies that $\{\mathbf{x}_i\}_{i=1}^n$ appear only in terms of their inner products, so we can non-linearize the algorithm by the *kernel trick* (Schölkopf & Smola, 2002), which is briefly explained below.

Let us consider a non-linear mapping $\phi(\mathbf{x})$ from \mathbb{R}^d to a reproducing kernel Hilbert space \mathcal{H} . Let $K(\mathbf{x}, \mathbf{x}')$ be the reproducing kernel of \mathcal{H} . A typical choice of the kernel function would be the Gaussian kernel:

$$K(\mathbf{x}, \mathbf{x}') = \exp(-\|\mathbf{x} - \mathbf{x}'\|^2/2\sigma^2). \quad (45)$$

For other choices, see, e.g., (Schölkopf & Smola, 2002). Because of the reproducing property of $K(\mathbf{x}, \mathbf{x}')$, \mathbf{K} is now the *kernel matrix*, i.e., the (i,j) -th element is defined as

$$\mathbf{K}_{i,j} = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle = K(\mathbf{x}_i, \mathbf{x}_j), \quad (46)$$

where $\langle \cdot, \cdot \rangle$ denotes the inner product in \mathcal{H} . Then the embedded image of $\phi(\mathbf{x}')$ by LFDA in \mathcal{H} is given by

$$(\bar{\boldsymbol{\alpha}}_1 | \bar{\boldsymbol{\alpha}}_2 | \cdots | \bar{\boldsymbol{\alpha}}_m)^\top \begin{pmatrix} K(\mathbf{x}_1, \mathbf{x}') \\ K(\mathbf{x}_2, \mathbf{x}') \\ \vdots \\ K(\mathbf{x}_n, \mathbf{x}') \end{pmatrix}, \quad (47)$$

where $\{\bar{\alpha}_i\}_{i=1}^d$ is the generalized eigenvectors associated to the generalized eigenvalues $\bar{\lambda}_1 \geq \bar{\lambda}_2 \geq \dots \geq \bar{\lambda}_d$ of the generalized eigenvalue problem (43). Note that $K\bar{L}^{(m)}K$ can be efficiently calculated by

$$K\bar{L}^{(m)}K = K\bar{D}^{(m)}K - (K\mathbf{1})(K\mathbf{1})^\top / n - K\bar{A}^{(m')}K. \quad (48)$$

4. Numerical Results

In this section, we apply the existing and proposed dimensionality reduction methods to benchmark data sets for visualization. In LPP and LFDA, the affinity matrix is determined by the same method used in Section 2.4.

We use *Iris*, *Letter recognition*, *Segment*, and *Thyroid disease* data sets available from the UCI machine learning repository. They are all multi-class classification data sets. We created two-class problems in the standard ‘one-versus-rest’ manner. Figure 2 shows the embedded samples in two-dimensional space by FDA, LFDA, and LPP.

For *Iris* data set, FDA tends to mix samples of different classes, while LFDA separates them very well. Multimodality of the ‘x’-class can be clearly observed in the results of LFDA, which implies that FDA did not work well because of multimodality. This experimental result supports the qualitative justification of LFDA given in Section 3.3. LPP also works well for this data set since three clusters (two ‘x’-class clusters and one ‘o’-class cluster) are well separated from each other in the original space.

For the other three data sets, FDA can separate samples of different classes quite well. LFDA also separates them equally well, but moreover it preserves multimodality of the ‘x’-class more prominently than FDA. This would be a desirable property in data visualization tasks. LPP also preserves multimodality of the data well, but it tends to mix samples of different classes.

Overall, LFDA is found to be useful in data visualization tasks.

5. Conclusions

Dimensionality reduction based on the Fisher criterion (FDA) works well given data samples of each class form a single cluster (i.e., unimodal). On the other hand, samples in some class can be multimodal, e.g., when multi-class classification problems are solved by a set of two-class ‘one-versus-rest’ problems. In this paper, we first showed that FDA can be regarded

as keeping in-class data pairs close and between-class data pairs apart. Based on this novel interpretation, we proposed a localized variant of FDA called local Fisher discriminant analysis (LFDA). We experimentally showed that LFDA can preserve multimodal structure of the data better than FDA, and LFDA can even provide more separate embedding than FDA.

Our future work includes the comparison with recently proposed methods, e.g., (Goldberger et al., 2005; Globerson & Roweis, 2006).

Acknowledgments

The author would like to thank Yuki Shinada, Klaus-Robert Müller, Hideki Asoh, Stefan Harmeling and anonymous reviewers for their comments. He also acknowledges financial support from MEXT (Grant-in-Aid for Young Scientists 17700142).

References

- Belkin, M., & Niyogi, P. (2003). Laplacian eigenmaps for dimensionality reduction and data representation. *Neural Computation*, 15, 1373–1396.
- Fisher, R. A. (1936). The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7, 179–188.
- Fukunaga, K. (1990). *Introduction to statistical pattern recognition*. Boston: Academic Press, Inc. Second edition.
- Globerson, A., & Roweis, S. (2006). Metric learning by collapsing classes. *Advances in Neural Information Processing Systems 18* (pp. 451–458). Cambridge, MA: MIT Press.
- Goldberger, J., Roweis, S., Hinton, G., & Salakhutdinov, R. (2005). Neighbourhood components analysis. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*, 513–520. Cambridge, MA: MIT Press.
- He, X., & Niyogi, P. (2004). Locality preserving projections. In S. Thrun, L. Saul and B. Schölkopf (Eds.), *Advances in neural information processing systems 16*. Cambridge, MA: MIT Press.
- Mika, S., Rätsch, G., Weston, J., Schölkopf, B., Smola, A., & Müller, K.-R. (2003). Constructing descriptive and discriminative nonlinear features: Rayleigh coefficients in kernel feature spaces. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25, 623–628.
- Schölkopf, B., & Smola, A. J. (2002). *Learning with kernels*. Cambridge, MA: MIT Press.
- Zelnik-Manor, L., & Perona, P. (2005). Self-tuning spectral clustering. In L. K. Saul, Y. Weiss and L. Bottou (Eds.), *Advances in neural information processing systems 17*, 1601–1608. Cambridge, MA: MIT Press.

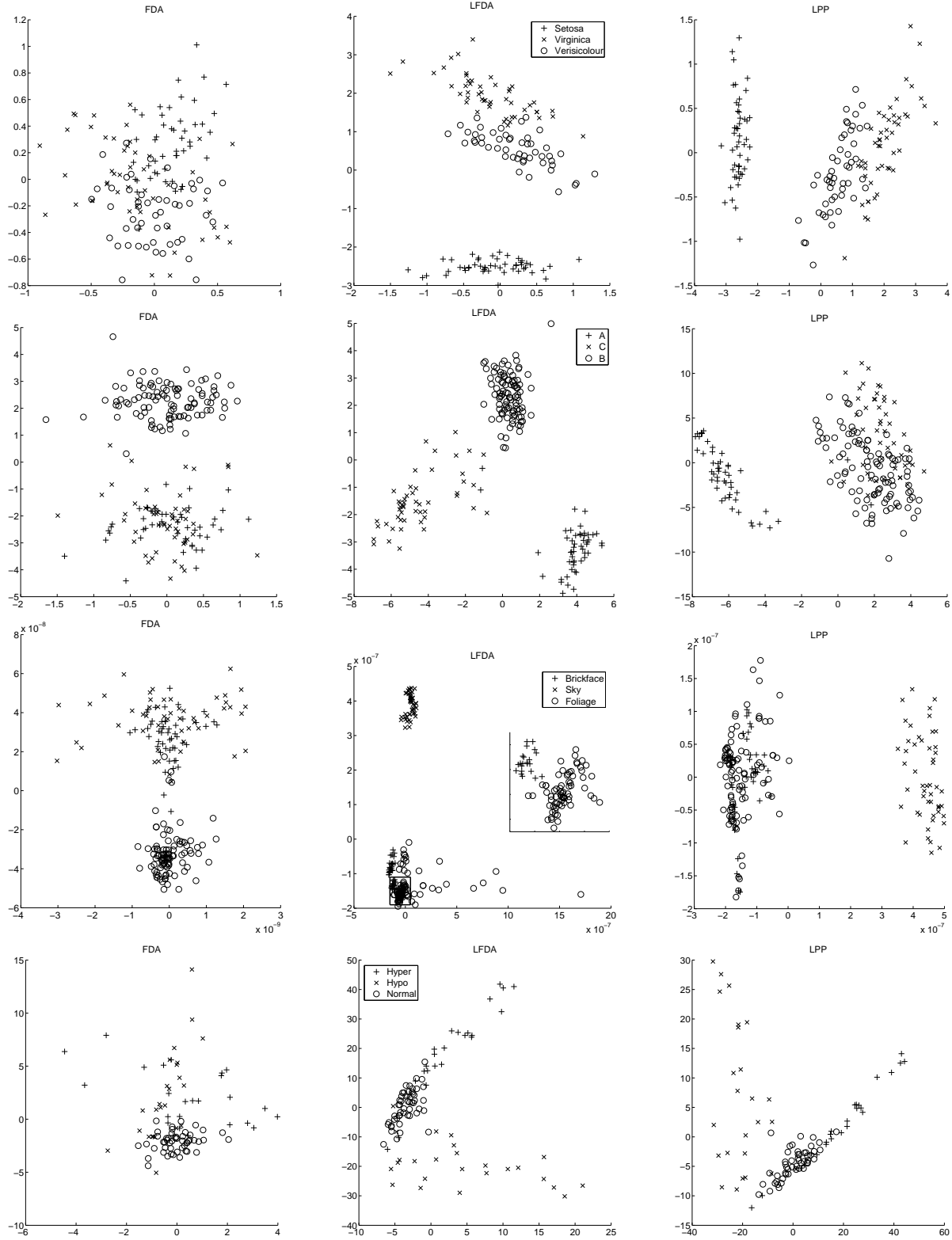


Figure 2. Results of data visualization ($m = 2$). From top to bottom, Iris data set ($d = 4$: ‘Setosa’ & ‘Virginica’ vs. ‘Versicolour’), Letter recognition data set ($d = 16$: ‘A’ & ‘C’ vs. ‘B’), Segment data set ($d = 18$: ‘Brickface’ & ‘Sky’ vs. ‘Foliage’), and Thyroid disease data set ($d = 5$: ‘Hyper’ & ‘Hypo’ vs. ‘Normal’). From left to right, FDA, LFDA, and LPP.