Advanced Data Analysis: Fisher Discriminant Analysis

Masashi Sugiyama (Computer Science)

W8E-505, <u>sugi@cs.titech.ac.jp</u> http://sugiyama-www.cs.titech.ac.jp/~sugi

## Supervised Dimensionality <sup>60</sup> Reduction

Samples  $\{x_i\}_{i=1}^n$  have class labels  $\{y_i\}_{i=1}^n$  :

 $\{(x_i, y_i)\}_{i=1}^n$ 

$$oldsymbol{x}_i \in \mathbb{R}^d \ y_i \in \{1, 2, \dots, c\}$$

We want to obtain an embedding such that samples in different classes are well separated from each other!





Setosa
Virginica
Verisicolour

#### Within-Class Scatter Matrix <sup>61</sup>

Sum of scatter within each class:

$$\boldsymbol{S}^{(w)} = \sum_{y=1}^{c} \sum_{i:y_i=y} (\boldsymbol{x}_i - \boldsymbol{\mu}_y) (\boldsymbol{x}_i - \boldsymbol{\mu}_y)^{\top}$$



$$\boldsymbol{\mu}_y = rac{1}{n_y} \sum_{i: y_i = y} \boldsymbol{x}_i$$

 $\mu_y$  :mean of samples in class y $n_y$  :# of samples in class y

### Between-Class Scatter Matrix <sup>62</sup>

#### Sum of scatter between classes:

$$\boldsymbol{S}^{(b)} = \sum_{y=1}^{c} n_y (\boldsymbol{\mu}_y - \boldsymbol{\mu}) (\boldsymbol{\mu}_y - \boldsymbol{\mu})^{\mathsf{T}}$$



$$oldsymbol{\mu} = rac{1}{n}\sum_{i=1}^n oldsymbol{x}_i$$

$$oldsymbol{\mu}_y = rac{1}{n_y} \sum_{i: y_i = y} oldsymbol{x}_i$$

 $\mu$  :mean of all samples  $\mu_y$  :mean of samples in class y $n_y$  :# of samples in class y

#### Fisher Discriminant Analysis (FDÅ) Fisher (1936)

- Idea: minimize within-class scatter and maximize between-class scatter by maximizing  $tr((BS^{(w)}B^{T})^{-1}BS^{(b)}B^{T})$
- To disable arbitrary scaling, we impose  $BS^{(w)}B^{\top} = I_m$

FDA criterion:

$$\boldsymbol{B}_{FDA} = \operatorname*{argmax}_{\boldsymbol{B} \in \mathbb{R}^{m \times d}} \operatorname{tr}(\boldsymbol{B}\boldsymbol{S}^{(b)}\boldsymbol{B}^{\top})$$

subject to  $\boldsymbol{B}\boldsymbol{S}^{(w)}\boldsymbol{B}^{\top} = \boldsymbol{I}_m$ 

## **FDA: Summary**

FDA criterion:  $B_{FDA} = \operatorname*{argmax}_{B \in \mathbb{R}^{m \times d}} \operatorname{tr}(BS^{(b)}B^{\top})$ subject to  $BS^{(w)}B^{\top} = I_m$ 

FDA solution:

$$\boldsymbol{B}_{FDA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^{\top}$$

•  $\{\lambda_i, \psi_i\}_{i=1}^m$ :Sorted generalized eigenvalues and normalized eigenvectors of  $S^{(b)}\psi = \lambda S^{(w)}\psi$ 

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \qquad \langle \boldsymbol{S}^{(w)} \boldsymbol{\psi}_i, \boldsymbol{\psi}_j \rangle = \delta_{i,j}$$

FDA embedding of a sample x:

$$z = B_{FDA}x$$

## **Examples of FDA** $d = 2, m = 1 \quad (\mathbb{R}^2 \Longrightarrow \mathbb{R}^1)$

65



FDA can find an appropriate subspace.



#### However, FDA does not work well if samples in a class have multimodality.

Dimensionality of Embedding Space We have  $\operatorname{rank}(S^{(b)}) = c - 1$ . (Homework) This means  $\{\lambda_i\}_{i=c}^d$  are always zero.  $\lambda_1 > \lambda_2 > \cdots > \lambda_d$ C :# of classes Due to the multiplicity of eigenvalues, eigenvectors  $\{\psi_i\}_{i=c}^d$  can be arbitrarily rotated in the null space of  $S^{(b)}$ . Thus FDA essentially requires m < c - 1

When c = 2, m can not be larger than 1 !

 $\ensuremath{\mathcal{M}}$  :dimensionality of embedding space

Pairwise Expressions of Scatter<sup>68</sup>  

$$S^{(w)} = \frac{1}{2} \sum_{i,j=1}^{n} Q_{i,j}^{(w)} (x_i - x_j) (x_i - x_j)^{\top}$$
 (Homework)  
 $Q_{i,j}^{(w)} = \begin{cases} \frac{1/n_y}{0} & (y_i = y_j = y) \\ 0 & (y_i \neq y_j) \end{cases}$   
 $S^{(b)} = \frac{1}{2} \sum_{i,j=1}^{n} Q_{i,j}^{(b)} (x_i - x_j) (x_i - x_j)^{\top}$   
 $Q_{i,j}^{(b)} = \begin{cases} \frac{1/n - 1/n_y}{0} & (y_i = y_j = y) \\ 1/n & (y_i \neq y_j) \end{cases}$ 

n :# of all samples

Implication:

 $n_y$  :# of samples in class y

- Samples in the same class are made close
- Samples in different classes are made apart

#### Local Fisher Discriminant Analysis Sugiyama (2007)

Idea: Take the locality of data into account:

- Nearby samples in the same class are made close
- Far-apart samples in the same class are not made close 10
- Samples in different classes are made apart



# **LFDA Criterion**

Local within-class scatter matrix:

$$\widetilde{\boldsymbol{S}}^{(w)} = \frac{1}{2} \sum_{i,j=1}^{n} \widetilde{\boldsymbol{Q}}_{i,j}^{(w)} (\boldsymbol{x}_i - \boldsymbol{x}_j) (\boldsymbol{x}_i - \boldsymbol{x}_j)^{\top} \qquad \boldsymbol{W}_{i,j} : \text{Similarity}$$
$$\widetilde{\boldsymbol{Q}}_{i,j}^{(w)} = \begin{cases} \boldsymbol{W}_{i,j} / n_y & (y_i = y_j = y) \\ 0 & (y_i \neq y_j) \end{cases}$$

Local between-class scatter matrix:

$$\widetilde{\boldsymbol{S}}^{(b)} = \frac{1}{2} \sum_{i,j=1}^{n} \widetilde{\boldsymbol{Q}}_{i,j}^{(b)} (\boldsymbol{x}_{i} - \boldsymbol{x}_{j}) (\boldsymbol{x}_{i} - \boldsymbol{x}_{j})^{\top}$$
$$\widetilde{\boldsymbol{Q}}_{i,j}^{(b)} = \begin{cases} \boldsymbol{W}_{i,j} (1/n - 1/n_{y}) & (y_{i} = y_{j} = y) \\ 1/n & (y_{i} \neq y_{j}) \end{cases}$$

LFDA criterion:  $B_{LFDA} = \operatorname*{argmax}_{B \in \mathbb{R}^{m \times d}} \operatorname{tr}(B\widetilde{S}^{(b)}B^{\top})$ 

subject to 
$$\boldsymbol{B}\widetilde{\boldsymbol{S}}^{(w)}\boldsymbol{B}^{\top} = \boldsymbol{I}_m$$

70

## LFDA: Summary

LFDA criterion:  $B_{LFDA} = \underset{B \in \mathbb{R}^{m \times d}}{\operatorname{argmax}} \operatorname{tr}(B\widetilde{S}^{(b)}B^{\top})$ subject to  $B\widetilde{S}^{(w)}B^{\top} = I_m$ 

LFDA solution:

$$\boldsymbol{B}_{LFDA} = (\boldsymbol{\psi}_1 | \boldsymbol{\psi}_2 | \cdots | \boldsymbol{\psi}_m)^\top$$

71

•  $\{\lambda_i, \psi_i\}_{i=1}^m$ : Sorted generalized eigenvalues and normalized eigenvectors of  $\widetilde{S}^{(b)}\psi = \lambda \widetilde{S}^{(w)}\psi$ 

$$\lambda_1 \geq \lambda_2 \geq \cdots \geq \lambda_d \qquad \langle \widetilde{m{S}}^{(w)} m{\psi}_i, m{\psi}_j 
angle = \delta_{i,j}$$

**LFDA** embedding of a sample x:

$$oldsymbol{z} = oldsymbol{B}_{LFDA}oldsymbol{x}$$



Note: Similarity matrix is defined by the nearestneighbor-based method with 50 nearest neighbors.

#### LFDA works well even for samples with within-class multimodality.

C :# of classes

Since  $\operatorname{rank}(\widetilde{\boldsymbol{S}}^{(b)}) \gg c$ , m can be large in LFDA.

 $\ensuremath{\mathcal{M}}$  :dimensionality of embedding space

# **Example of FDA/LFDA**

- Thyroid disease data (5-dimensional)
  - T3-resin uptake test.
  - Total Serum thyroxin as measured by the isotopic displacement method.

etc

- Label: Healty or sick
- Sick can caused by
  - Hyper-functioning of thyroid
  - Hypo-functioning of thyroid

#### Projected Samples onto 1-D Space FDA LFDA

Sick









- Sick and healthy are nicely split.
- But hyper- and hypofunctioning are mixed.

- Sick and healthy are nicely split.
- Hyper- and hypo-functioning are also nicely separated.

#### Homework

#### 1. Implement FDA/LFDA and reproduce the 2dimensional examples shown in the class.

http://sugiyama-www.cs.titech.ac.jp/~sugi/data/DataAnalysis



Test FDA/LFDA with your own (artificial or real) data and analyze the characteristics of FDA/LFDA.

# Homework (cont.)

76

2. Prove that  $\operatorname{rank}(S^{(b)}) = c - 1$ . *c* :# of classes Hint: Range of  $S^{(b)}$  is spanned by  $\{\mu_y - \mu\}_{y=1}^c$ .



# Homework (cont.)

//

3. Prove that A)  $S^{(w)} = rac{1}{2} \sum_{i=1}^{n} Q_{i,j}^{(w)} (x_i - x_j) (x_i - x_j)^{ op}$ B)  $S^{(b)} = \frac{1}{2} \sum_{i=1}^{n} Q_{i,j}^{(b)} (x_i - x_j) (x_i - x_j)^{\top}$  $\boldsymbol{Q}_{i,j}^{(w)} = \begin{cases} 1/n_y & (y_i = y_j = y) \\ 0 & (y_i \neq y_j) \end{cases} \boldsymbol{Q}_{i,j}^{(b)} = \begin{cases} 1/n - 1/n_y & (y_i = y_j = y) \\ 1/n & (y_i \neq y_j) \end{cases}$  $n_y$  :# of samples in class y = n :# of all samples

Hint: The use of the following mixture scatter matrix may make your life easy...

$$oldsymbol{S}^{(m)} = oldsymbol{S}^{(w)} + oldsymbol{S}^{(b)} \left( = \sum_{i=1}^n (oldsymbol{x}_i - oldsymbol{\mu}) (oldsymbol{x}_i - oldsymbol{\mu})^ op 
ight)$$

### Suggestion

Read the following article for the next class:

 B. Schölkopf, A. Smola and K.-R. Müller: Nonlinear Component Analysis as a Kernel Eigenvalue Problem, *Neural Computation*, 10(5), 1299-1319, 1998.

http://neco.mitpress.org/cgi/reprint/10/5/1299.pdf