

EE5907: Pattern Recognition & Machine Learning EE5026: for Data Analytics

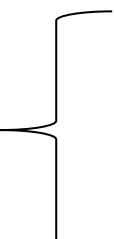
Dr. WANG Si

Email: si.wang@nus.edu.sg

Office: E1a-04-01

A Gentle Reminder

CA2 for EE5907 and CA1 for EE5026 are released!

- Face Recognition System (CMU PIE dataset) 
 - Part 1: PCA & LDA (from scratch)
 - Part 2: KNN, GMM, SVM, CNN
- Go through Assignment with the instruction file.
- Deadline (strict): 14 Nov. 2025 1800 SGT
- Individual assignment

Outlines

- Pattern Representation Learning
 - Unsupervised Representation Learning (week 7)
 - Supervised Representation Learning (week 8)
 - Unified Framework: Graph Embedding (week 8)
- Pattern Recognition Models
 - Clustering and Applications (week 9)
 - Gaussian Mixture Model and Boosting (week 10)
 - Support Vector Machines (week 11)
- Deep Learning (week 12)
- Revision and Q&A (week 13)

Pattern Representation Learning

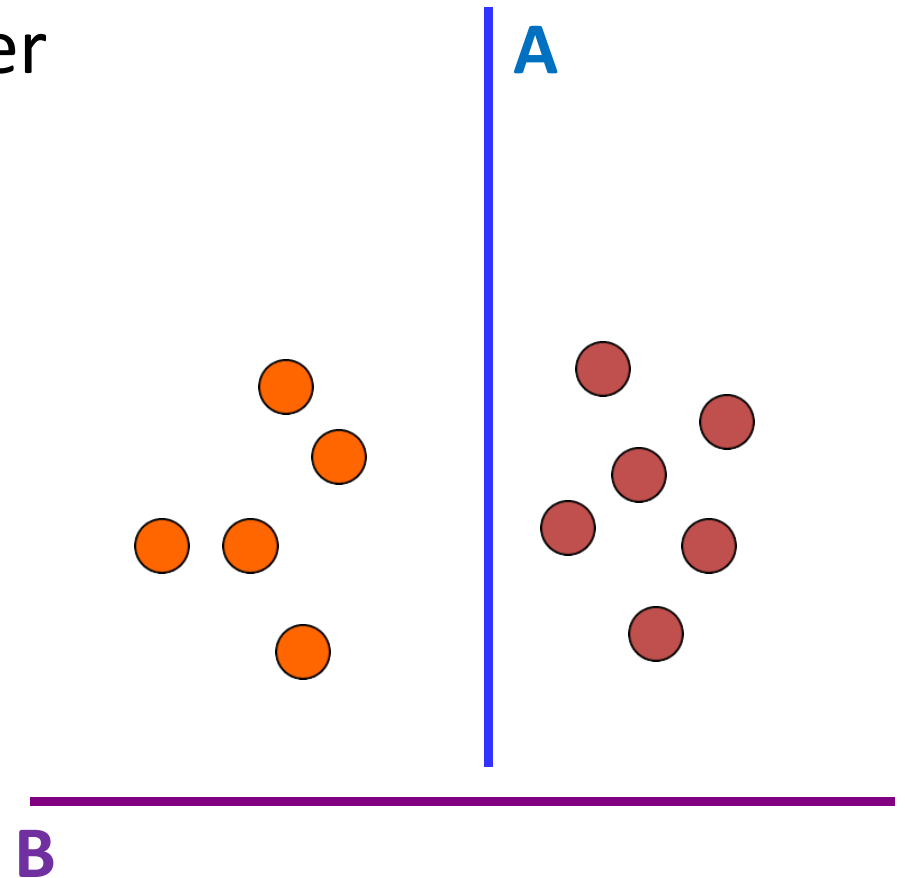
- We will discuss:
 - Linear Discriminant Analysis (LDA)
 - A Unified Framework: Graph Embedding (GE)
 - An Example of GE: Marginal Fisher Analysis (MFA)
- At the end of this lecture, you should be able to:
 - Understand rationales behind LDA and GE
 - Perform LDA

Supervised Feature Learning

Linear Discriminant Analysis

What is a Good Projection?

- What is a good criterion under supervised setting?
 - Separating different classes

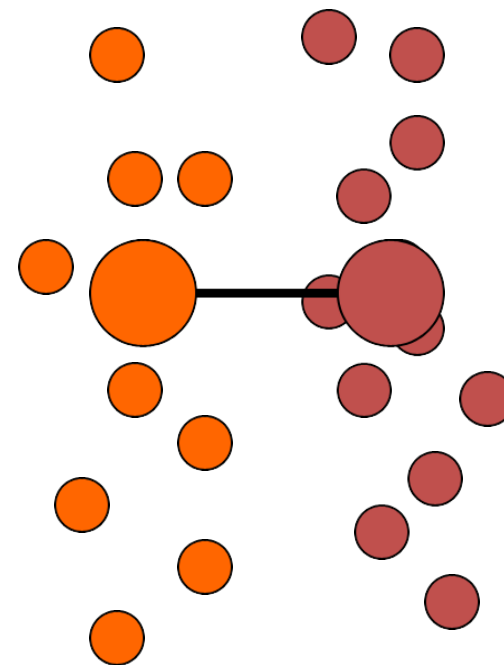


Which projection axis is better under supervised setting?

What Class Information May be Useful?

- Between-class distance
 - Distance between the centroids of different classes

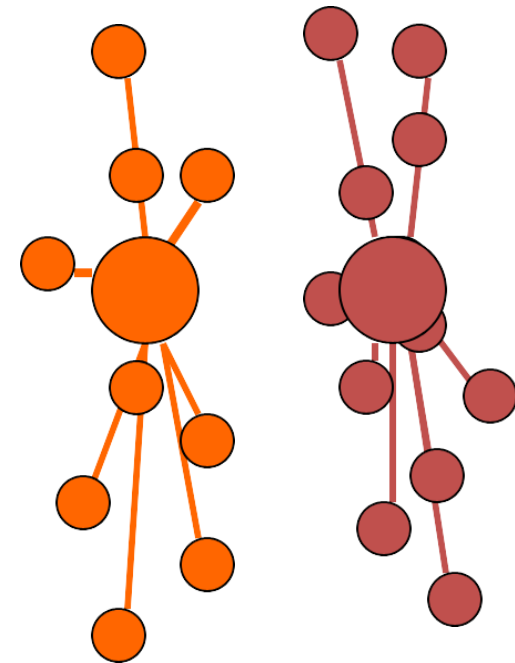
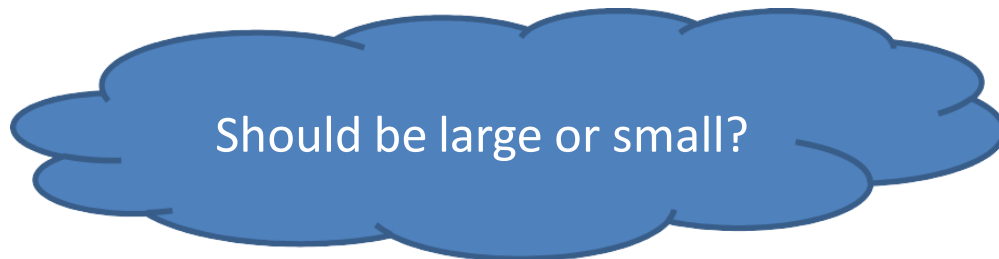
Should be large or small?



— Between-class distance

What Class Information May be Useful?

- Between-class distance
 - Distance between the centroids of different classes
- Within-class distance
 - Accumulated distance of an instance to the centroid of its class

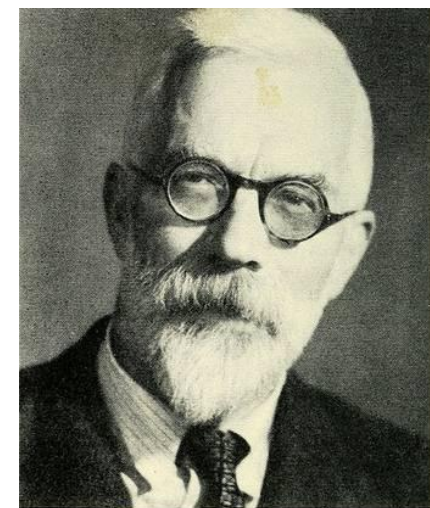


— Within-class distance

Linear Discriminant Analysis!

Linear Discriminant Analysis

- A bit history
 - Fisher's Linear Discriminative Analysis (1936)
 - Fisher, R. A. (1936). "The Use of Multiple Measurements in Taxonomic Problems". Annals of Eugenics.
 - Linear Discriminative Analysis (1948)
 - Rao, R. C. (1948). "The utilization of multiple measurements in problems of biological classification". Journal of the Royal Statistical Society, Series B.

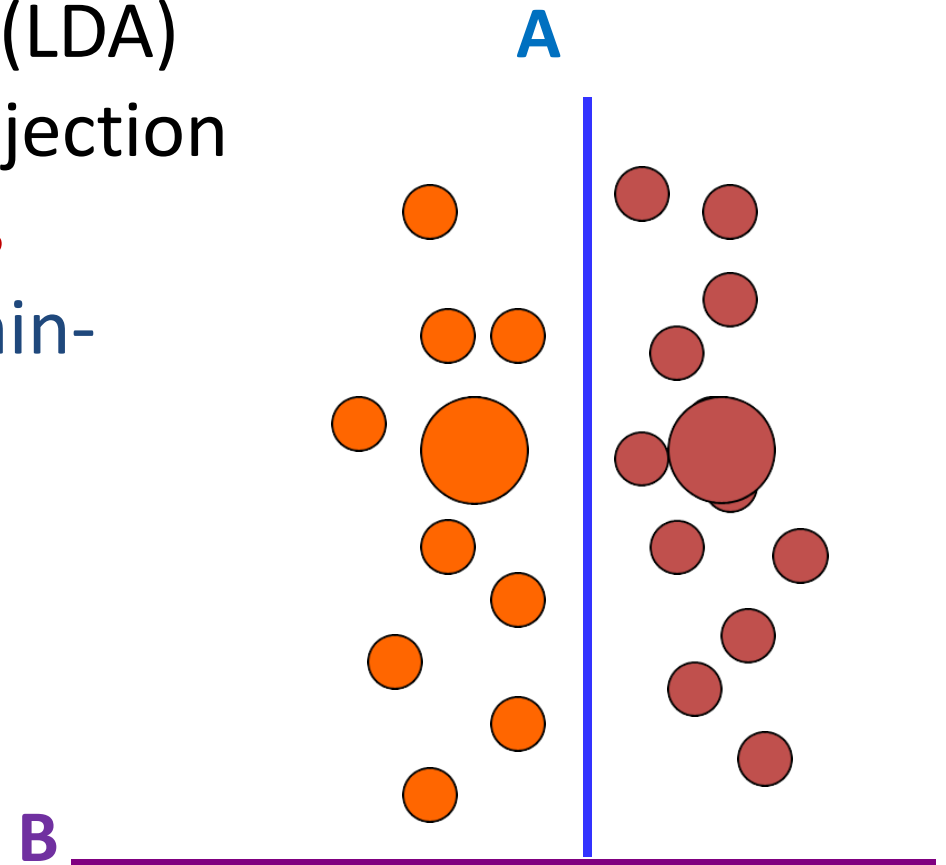


Ronald Fisher (1890-1962)

"A genius who almost single-handedly created the foundations for **modern statistical science**" – Anders Hald

Linear Discriminant Analysis

- **Linear discriminant analysis** (LDA) finds most discriminative projection by **maximizing between-class distance** and **minimizing within-class distance**



Which projection axis is possibly the 1st PC, which one is possibly the LDA projection?

Statistical Facts (1)

Quantity for Measuring Within and Between Class Distance

- Total mean vector:

$$\mu = \frac{1}{N} \sum_{\mathbf{x}} \mathbf{x}, \text{ } N \text{ is the number of all samples.}$$

- Class-specific mean vector:

$$\mu_i = \frac{1}{n_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}, \text{ } n_i \text{ is the size of class } C_i.$$

- Class-specific scatter matrix $\hat{\mathbf{S}}_i$:

$$\hat{\mathbf{S}}_i = \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T$$

Statistical Facts (2)

Quantity for Measuring Within and Between Class Distance (cont.)

- Within-class scatter matrix $\widehat{\mathcal{S}}_W$:

$$\widehat{\mathcal{S}}_W = \sum_{i=1}^C \widehat{\mathcal{S}}_i \quad C \text{ is the class number.}$$

- Between-class scatter matrix $\widehat{\mathcal{S}}_B$:

$$\widehat{\mathcal{S}}_B = \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

- Total scatter matrix $\widehat{\mathcal{S}}_T$:

$$\begin{aligned} \widehat{\mathcal{S}}_T &= \sum_x (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \\ &= \widehat{\mathcal{S}}_W + \widehat{\mathcal{S}}_B \end{aligned}$$



see next slide

Proof

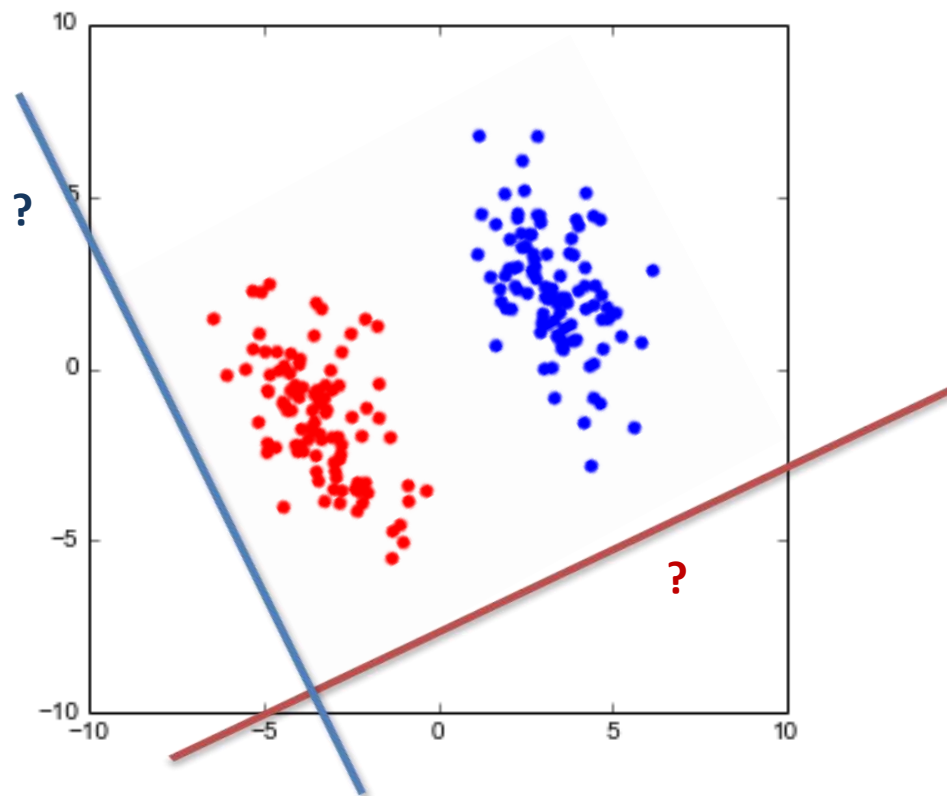
$$\widehat{\mathbf{S}}_T = \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T \quad \widehat{\mathbf{S}}_w = \sum_{i=1}^C \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T \quad \widehat{\mathbf{S}}_B = \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

$$\begin{aligned} \widehat{\mathbf{S}}_w + \widehat{\mathbf{S}}_B &= \sum_{i=1}^C \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T + \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T \\ &= \sum_{i=1}^C (\sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T + n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T) \\ &= \sum_{i=1}^C (\sum_{\mathbf{x} \in C_i} \mathbf{x}\mathbf{x}^T - 2\mathbf{x}\boldsymbol{\mu}_i^T + \boldsymbol{\mu}_i\boldsymbol{\mu}_i^T + n_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T - 2n_i\boldsymbol{\mu}_i\boldsymbol{\mu}^T + n_i\boldsymbol{\mu}\boldsymbol{\mu}^T) \\ &= \sum_{i=1}^C (n_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T - 2n_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T + \sum_{\mathbf{x} \in C_i} \mathbf{x}\mathbf{x}^T + n_i\boldsymbol{\mu}_i\boldsymbol{\mu}_i^T - 2n_i\boldsymbol{\mu}_i\boldsymbol{\mu}^T + n_i\boldsymbol{\mu}\boldsymbol{\mu}^T) \\ &= \sum_{i=1}^C (\sum_{\mathbf{x} \in C_i} \mathbf{x}\mathbf{x}^T - 2\sum_{\mathbf{x} \in C_i} \mathbf{x}\boldsymbol{\mu}_i^T + n_i\boldsymbol{\mu}\boldsymbol{\mu}^T) \\ &= \sum_{\mathbf{x}} \mathbf{x}\mathbf{x}^T - 2\mathbf{x}\boldsymbol{\mu}^T + \boldsymbol{\mu}\boldsymbol{\mu}^T \\ &= \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \mathbf{S}_T \end{aligned}$$

within-class scatter + between-class scatter
= total scatter.

Linear Discriminant Analysis - Problem

- **PROBLEM: To separate populations**
 - Suppose we have C classes from D -dimensional distributions. We want to project these classes onto a p -dimensional subspace ($p < D$) so that the variation **between the classes** is **as large as possible**, relative to the **variation within the classes**.



LDA: Two-class

A simple case

- Assume we have a set of d -dimensional samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, n_1 of which from class C_1 , and n_2 of which from class C_2
- We seek for a linear projection vector \mathbf{w} , along which we project the samples \mathbf{x} onto a line and get a scalar y :

$$y = \mathbf{w}^T \mathbf{x}$$

such that, the separability of the projected samples between 2 classes are maximized.

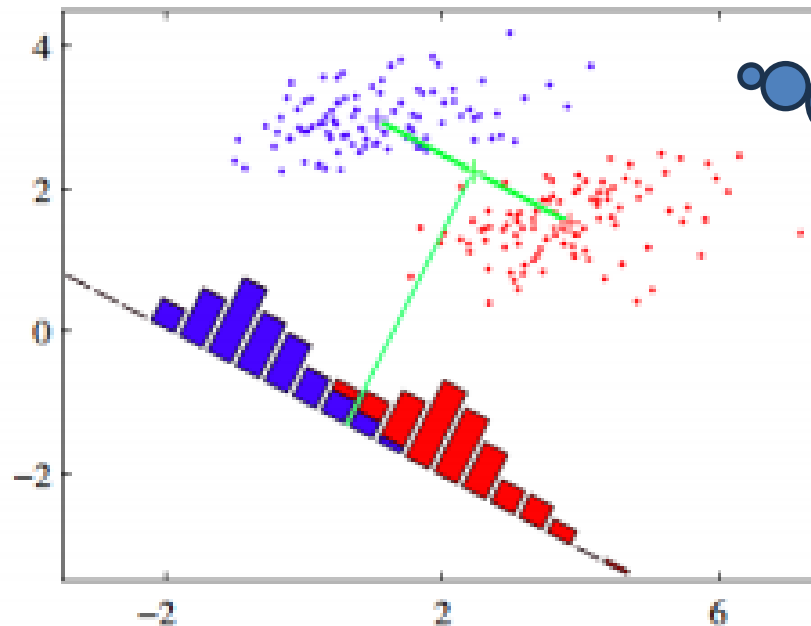
LDA: Two-class

- If the class separability measure is only **maximizing the distance between projected means** $\tilde{\mu}_1$ and $\tilde{\mu}_2$:

$$J(\mathbf{w}) = |\tilde{\mu}_1 - \tilde{\mu}_2| = |\mathbf{w}^T(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)|$$

where $\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}$ and $\tilde{\mu}_i = \mathbf{w}^T \boldsymbol{\mu}_i$

$1/|\mathbf{w}| \cdot \cos \theta$
 $\sim 1/|\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2|$



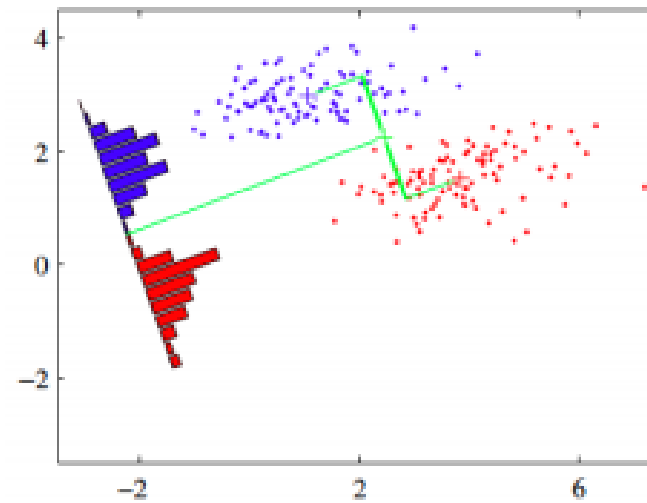
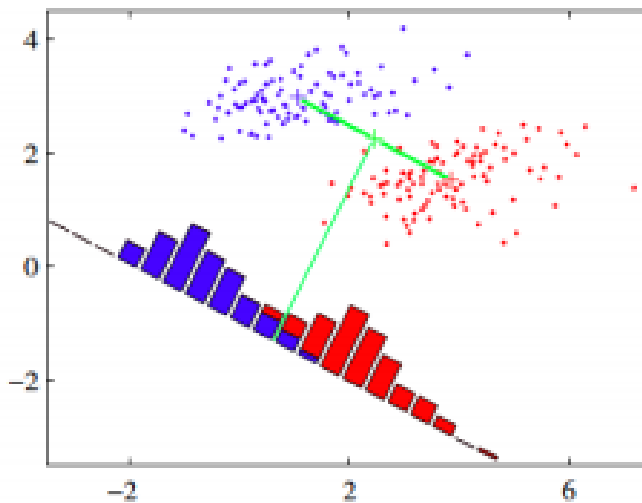
Is it the optimal solution?

LDA: Two-class

- If the class separability measure is only **maximizing the distance between projected means** $\tilde{\mu}_1$ and $\tilde{\mu}_2$:

$$J(\mathbf{w}) = |\tilde{\mu}_1 - \tilde{\mu}_2| = |\mathbf{w}^T(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)|$$

where $\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{x \in C_i} \mathbf{x}$ and $\tilde{\mu}_i = \mathbf{w}^T \boldsymbol{\mu}_i$



LDA: Two-class

- Fisher suggested **maximizing the difference between the projected means**, normalized by **a measure of the projected within-class scatter!**

- For each class, the projected scatter is:

$$\tilde{S}_i = \sum_{y \in C_i} (y - \tilde{\mu}_i)^2$$

class-specific scatter
 $\sum_{x \in C_i} (x - \tilde{\mu}_i)^2$ max. diff.

then, $(\tilde{S}_1 + \tilde{S}_2)$ is the within-class scatter of the projected samples.

- Therefore, the criterion function: $J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{S}_1 + \tilde{S}_2}$

$$\text{rank}(X) = \text{rank}(XX^T) = \text{rank}(X^T X)$$

LDA: Two-class

- To find the optimum \mathbf{w}^* , we need write $J(\mathbf{w})$ as a function of \mathbf{w} :

$$J(\mathbf{w}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{S}_1 + \tilde{S}_2} = \frac{(\mathbf{w}^T(\mu_1 - \mu_2))^2}{\sum_{y \in C_1} (y - \tilde{\mu}_1)^2 + \sum_{y \in C_2} (y - \tilde{\mu}_2)^2}$$

$$= \frac{\mathbf{w}^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \mathbf{w}}{\sum_{i=1}^C \sum_{x \in C_i} (\mathbf{w}^T \mathbf{x} - \mathbf{w}^T \mu_i)^2}$$

$$= \frac{\mathbf{w}^T \left[\frac{n_1 + n_2}{n_1 n_2} \widehat{\mathbf{S}}_B \right] \mathbf{w}}{\mathbf{w}^T \left[\sum_{i=1}^C \sum_{x \in C_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T \right] \mathbf{w}}$$

$$= \frac{n_1 + n_2}{n_1 n_2} \frac{\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w}}{\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w}}$$

$$\sum_i h_i (\underbrace{\mathbf{w}^T \mu_i}_{\text{✓}} - \underbrace{\mathbf{w}^T \mu_i}_{\text{✓}})^2$$

see next slide

$$\mathbf{w}^T \sum_i n_i (\mathbf{x}_i - \mu_i)(\mathbf{x}_i - \mu_i)^T \mathbf{w}$$

Note that $\widehat{\mathbf{S}}_B$ has a rank of at most 1.

Proof

• Prove $\widehat{\mathbf{S}}_B = \frac{n_1 n_2}{n_1 + n_2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$

$$\widehat{\mathbf{S}}_B = \sum_{i=1}^c n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

$$= n_1 (\boldsymbol{\mu}_1 - \boldsymbol{\mu})(\boldsymbol{\mu}_1 - \boldsymbol{\mu})^T + n_2 (\boldsymbol{\mu}_2 - \boldsymbol{\mu})(\boldsymbol{\mu}_2 - \boldsymbol{\mu})^T$$

$$= n_1 (\boldsymbol{\mu}_1 \boldsymbol{\mu}_1^T - 2\boldsymbol{\mu}_1 \boldsymbol{\mu}^T + \boldsymbol{\mu} \boldsymbol{\mu}^T) + n_2 (\boldsymbol{\mu}_2 \boldsymbol{\mu}_2^T - 2\boldsymbol{\mu}_2 \boldsymbol{\mu}^T + \boldsymbol{\mu} \boldsymbol{\mu}^T)$$

$$(n_1 + n_2) \boldsymbol{\mu} = n_1 \boldsymbol{\mu}_1 + n_2 \boldsymbol{\mu}_2$$

$$= n_1 \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^T + n_2 \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^T - (n_1 + n_2) \boldsymbol{\mu} \boldsymbol{\mu}^T$$

$$= n_1 \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^T + n_2 \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^T - \frac{(n_1 \boldsymbol{\mu}_1 + n_2 \boldsymbol{\mu}_2)(n_1 \boldsymbol{\mu}_1 + n_2 \boldsymbol{\mu}_2)^T}{n_1 + n_2}$$

$$= \frac{n_1 n_2 \boldsymbol{\mu}_1 \boldsymbol{\mu}_1^T + n_1 n_2 \boldsymbol{\mu}_2 \boldsymbol{\mu}_2^T - 2n_1 n_2 \boldsymbol{\mu}_1 \boldsymbol{\mu}_2^T}{n_1 + n_2}$$

$$= \frac{n_1 n_2}{n_1 + n_2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$$

Handwritten notes:

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T = \begin{bmatrix} \mu_{11} - \mu_{21} & \mu_{12} - \mu_{22} \\ \vdots & \vdots \end{bmatrix} \begin{bmatrix} \mu_{11} - \mu_{21} & \mu_{12} - \mu_{22} & \dots & \mu_{1d} - \mu_{2d} \\ \mu_{21} - \mu_{11} & \mu_{22} - \mu_{12} & \dots & \mu_{2d} - \mu_{1d} \\ \vdots & \vdots & \ddots & \vdots \\ \mu_{1d} - \mu_{2d} & \mu_{2d} - \mu_{1d} & \dots & \mu_{dd} - \mu_{dd} \end{bmatrix}$$

Example for $d=2$:

$$= \begin{bmatrix} \mu_{11} - \mu_{21} & \mu_{12} - \mu_{22} \\ \mu_{21} - \mu_{11} & \mu_{22} - \mu_{12} \\ \mu_{12} - \mu_{22} & \mu_{22} - \mu_{12} \\ \mu_{22} - \mu_{12} & \mu_{12} - \mu_{22} \end{bmatrix}$$

LDA: Two-class

- We can finally express the **Fisher criterion** in terms of $\widehat{\mathbf{S}}_W$ and $\widehat{\mathbf{S}}_B$ as

$$J(\mathbf{w}) = \frac{\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w}}{\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w}}$$

Ignore $\frac{n_1+n_2}{n_1 n_2}$

- To find the maximum of $J(\mathbf{w})$, we take derivative and set to zero

$$\frac{d}{d\mathbf{w}} J(\mathbf{w}) = 0$$

$$\frac{d}{dx} \left(\frac{f(x)}{g(x)} \right) = \frac{g(x) \frac{d}{dx} f(x) - f(x) \frac{d}{dx} g(x)}{(g(x))^2}$$

$$\frac{(\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w}) \frac{d}{d\mathbf{w}} (\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w}) - (\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w}) \frac{d}{d\mathbf{w}} (\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w})}{(\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w})^2} = 0$$

$$2(\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w}) \widehat{\mathbf{S}}_B \mathbf{w} - 2(\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w}) \widehat{\mathbf{S}}_W \mathbf{w} = 0$$

Divide by
 $(\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w})$

$$\widehat{\mathbf{S}}_B \mathbf{w} - \frac{(\mathbf{w}^T \widehat{\mathbf{S}}_B \mathbf{w})}{(\mathbf{w}^T \widehat{\mathbf{S}}_W \mathbf{w})} \widehat{\mathbf{S}}_W \mathbf{w} = 0$$

$$\widehat{\mathbf{S}}_B \mathbf{w} - J(\mathbf{w}) \widehat{\mathbf{S}}_W \mathbf{w} = 0 \Rightarrow \widehat{\mathbf{S}}_B \mathbf{w} = J(\mathbf{w}) \widehat{\mathbf{S}}_W \mathbf{w}$$

$$\frac{d\mathbf{w}^T A \mathbf{w}}{d\mathbf{w}} = (A + A^T) \mathbf{w}$$

LDA: Two-class

- The **generalized eigenvalue** problem ($A\mathbf{v} = \lambda B\mathbf{v}$)

$$\widehat{\mathbf{S}}_B \mathbf{w} = J(\mathbf{w}) \widehat{\mathbf{S}}_W \mathbf{w}$$

- If $\widehat{\mathbf{S}}_W$ is invertible, convert to standard eigenvalue problem ($A\mathbf{v} = \lambda \mathbf{v}$):

$$\widehat{\mathbf{S}}_W^{-1} \widehat{\mathbf{S}}_B \mathbf{w} = J(\mathbf{w}) \mathbf{w}$$

Then, \mathbf{w}^* is the eigenvector of $\widehat{\mathbf{S}}_W^{-1} \widehat{\mathbf{S}}_B$.

Its corresponding eigenvalue, i.e., $J(\mathbf{w})^*$, should be the largest.

- If $\widehat{\mathbf{S}}_W$ is not invertible,

$$(\widehat{\mathbf{S}}_B - J(\mathbf{w}) \widehat{\mathbf{S}}_W) \mathbf{w} = \mathbf{0}$$

$$\det(\widehat{\mathbf{S}}_B - J(\mathbf{w}) \widehat{\mathbf{S}}_W) = 0$$

Examples: LDA: Two-class

Given the following dataset, compute the LDA projection vector \mathbf{w}

Sample	Feature 1	Feature 2	Class
1	2	3	1
2	3	3	1
3	2	4	1
4	6	7	2
5	7	7	2
6	6	8	2

Examples: LDA: Two-class

Step 1: Compute μ, μ_1, μ_2

$$\mu = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i = \frac{1}{6} \left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 3 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} + \begin{bmatrix} 6 \\ 7 \end{bmatrix} + \begin{bmatrix} 7 \\ 7 \end{bmatrix} + \begin{bmatrix} 6 \\ 8 \end{bmatrix} \right) = \begin{bmatrix} 4.33 \\ 5.33 \end{bmatrix}$$

$$\mu_1 = \frac{1}{3} \sum_{i=1}^3 \mathbf{x}_i = \frac{1}{3} \left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} + \begin{bmatrix} 3 \\ 3 \end{bmatrix} + \begin{bmatrix} 2 \\ 4 \end{bmatrix} \right) = \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix}$$

$$\mu_2 = \frac{1}{3} \sum_{i=4}^6 \mathbf{x}_i = \frac{1}{3} \left(\begin{bmatrix} 6 \\ 7 \end{bmatrix} + \begin{bmatrix} 7 \\ 7 \end{bmatrix} + \begin{bmatrix} 6 \\ 8 \end{bmatrix} \right) = \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix}$$

N=6, n₁=3, n₂=3

Sample	Feature 1	Feature 2	Class
1	2	3	1
2	3	3	1
3	2	4	1
4	6	7	2
5	7	7	2
6	6	8	2

Step 2: Compute $\widehat{\mathcal{S}}_W, \widehat{\mathcal{S}}_B$

$$\begin{aligned} \widehat{\mathcal{S}}_W &= \sum_{i=1}^C \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)^T = \left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right) \left(\begin{bmatrix} 2 \\ 3 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right)^T + \left(\begin{bmatrix} 3 \\ 3 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right) \left(\begin{bmatrix} 3 \\ 3 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right)^T \\ &+ \left(\begin{bmatrix} 2 \\ 4 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right) \left(\begin{bmatrix} 2 \\ 4 \end{bmatrix} - \begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} \right)^T + \left(\begin{bmatrix} 6 \\ 7 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right) \left(\begin{bmatrix} 6 \\ 7 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right)^T + \left(\begin{bmatrix} 7 \\ 7 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right) \left(\begin{bmatrix} 7 \\ 7 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right)^T \\ &+ \left(\begin{bmatrix} 6 \\ 8 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right) \left(\begin{bmatrix} 6 \\ 8 \end{bmatrix} - \begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} \right)^T = \begin{bmatrix} 1.33 & -0.67 \\ -0.67 & 1.33 \end{bmatrix} \end{aligned}$$

$$\begin{aligned} \widehat{\mathcal{S}}_B &= \sum_{i=1}^C n_i (\mu_i - \mu)(\mu_i - \mu)^T \\ &= 3 \left(\begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} - \begin{bmatrix} 4.33 \\ 5.33 \end{bmatrix} \right) \left(\begin{bmatrix} 2.33 \\ 3.33 \end{bmatrix} - \begin{bmatrix} 4.33 \\ 5.33 \end{bmatrix} \right)^T + 3 \left(\begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} - \begin{bmatrix} 4.33 \\ 5.33 \end{bmatrix} \right) \left(\begin{bmatrix} 6.33 \\ 7.33 \end{bmatrix} - \begin{bmatrix} 4.33 \\ 5.33 \end{bmatrix} \right)^T = \begin{bmatrix} 24 & 24 \\ 24 & 24 \end{bmatrix} \end{aligned}$$

Examples: LDA: Two-class

Step 3: Compute eigenvector \mathbf{w}^*

$$(\widehat{\mathbf{S}}_B - J(\mathbf{w})\widehat{\mathbf{S}}_W)\mathbf{w} = \mathbf{0}$$

$$\det(\widehat{\mathbf{S}}_B - J(\mathbf{w})\widehat{\mathbf{S}}_W) = 0$$

$$\det\left(\begin{bmatrix} 24 & 24 \\ 24 & 24 \end{bmatrix} - J(\mathbf{w})\begin{bmatrix} 1.33 & -0.67 \\ -0.67 & 1.33 \end{bmatrix}\right) = 0$$

$$\det\left(\begin{bmatrix} 24 - 1.33J(\mathbf{w}) & 24 + 0.67J(\mathbf{w}) \\ 24 + 0.67J(\mathbf{w}) & 24 - 1.33J(\mathbf{w}) \end{bmatrix}\right) = 0$$

$$(24 - 1.33J(\mathbf{w}))^2 - (24 + 0.67J(\mathbf{w}))^2 = 0$$

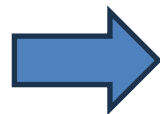
$$J_1 = 72, \quad J_2 = 0$$

Since $J(\mathbf{w})$ should be maximized, $J(\mathbf{w}) = J_1 = 72$

$$\text{Then, } (\widehat{\mathbf{S}}_B - J(\mathbf{w})\widehat{\mathbf{S}}_W)\mathbf{w} = \mathbf{0}$$

$$\begin{bmatrix} -72 & 72 \\ 72 & -72 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \mathbf{0}$$

$$w_1 = w_2$$



After normalization, $\mathbf{w}^* = \pm \begin{bmatrix} \frac{1}{\sqrt{2}} \\ \frac{1}{\sqrt{2}} \end{bmatrix}$

Sample	Feature 1	Feature 2	Class
1	2	3	1
2	3	3	1
3	2	4	1
4	6	7	2
5	7	7	2
6	6	8	2

LDA: Multi-class

- Instead of **one projection** y , we will now seek at most $p = \min(d, C - 1)$ projections $[y_1, y_2, \dots, y_p]$ by means of p projection vectors \mathbf{w}_i arranged by columns into a projection matrix $\mathbf{W} = [\mathbf{w}_1 | \mathbf{w}_2 | \dots | \mathbf{w}_p]$:

$$y_i = \mathbf{w}_i^T \mathbf{x} \quad \Rightarrow \quad \mathbf{y} = \mathbf{W}^T \mathbf{x}$$

- From our derivation for the two-class problem, we have the scatter matrices for the projected samples

$$\tilde{\mathbf{S}}_W = \mathbf{W}^T \widehat{\mathbf{S}}_W \mathbf{W}$$

$$\tilde{\mathbf{S}}_B = \mathbf{W}^T \widehat{\mathbf{S}}_B \mathbf{W}$$

LDA: Multi-class

- we use the **determinant** of the scatter matrices to obtain a scalar objective function

$$J(W) = \frac{\det(\tilde{S}_B)}{\det(\tilde{S}_W)} = \frac{\det(W^T \hat{S}_B W)}{\det(W^T \hat{S}_W W)}$$

- The optimal projection matrix W^* that maximizes $J(W)$ is the one whose columns are the **eigenvectors** w_i^* corresponding to the **largest eigenvalues** of $\hat{S}_W^{-1} \hat{S}_B$ (if \hat{S}_W is invertible)

$$W^* = \operatorname{argmax} \frac{\det(W^T \hat{S}_B W)}{\det(W^T \hat{S}_W W)} \rightarrow \boxed{\hat{S}_W^{-1} \hat{S}_B w_i^* = J(w_i^*) w_i^*}$$

OR

$$\boxed{\hat{S}_B w_i^* = J(w_i^*) \hat{S}_W w_i^*}$$

Classification with LDA

- Compute LDA projection vectors using training dataset
- Project both training and testing data into the obtained LDA space
- For each testing data point, use K-Nearest-Neighbor (KNN) or other classifiers for classification

Basis Visualization



Sample Images from YALE database



(a)



(b)

The first 10 (a) Eigenfaces, (b) Fisherfaces, calculated from the face images in the YALE database.

↓
PCA

↓
LDA

Some Remarks

- Unlike principal component analysis (PCA), the linear discriminant transformation W is **not necessarily orthogonal**.
 - Because $\widehat{\mathbf{S}}_W^{-1}\widehat{\mathbf{S}}_B$ may not be symmetric.
- Why is the number of the LDA projections at most $\min(d, C-1)$?

Some Remarks

- Why is the number of the LDA projections at most $\min(d, C-1)$?



meaningful generalized eigenvectors of $(\widehat{\mathbf{S}}_B, \widehat{\mathbf{S}}_W)$ is at most $\min(d, C-1)$
(Nonzero eigenvalues)

$$\widehat{\mathbf{S}}_B \mathbf{w}_i^* = J(\mathbf{w}_i^*) \widehat{\mathbf{S}}_W \mathbf{w}_i^*$$

meaningful generalized eigenvectors is limited by $\text{Rank}(\widehat{\mathbf{S}}_B)$

$$\widehat{\mathbf{S}}_B = \sum_{i=1}^C n_i \underbrace{(\mu_i - \mu)(\mu_i - \mu)^T}_{\text{rank} \leq 1}$$

Rank subadditivity rule

$$\text{rank}(A + B) \leq \text{rank}(A) + \text{rank}(B)$$

$$\text{rank}(\widehat{\mathbf{S}}_B) \leq C$$

Rank can't exceed
the matrix dimension
 d

$$\sum_{i=1}^C n_i (\mu_i - \mu) = \mathbf{0} \Rightarrow \sum_{i=1}^C n_i \mu_i = \mu \sum_{i=1}^C n_i \Rightarrow \sum_{i=1}^C n_i \mu_i = \bar{\mu} \sum_{i=1}^C n_i$$

$$\text{Rank}(\widehat{\mathbf{S}}_B) \leq C-1$$

$$\text{Rank}(\widehat{\mathbf{S}}_B) \leq \min(d, C-1)$$

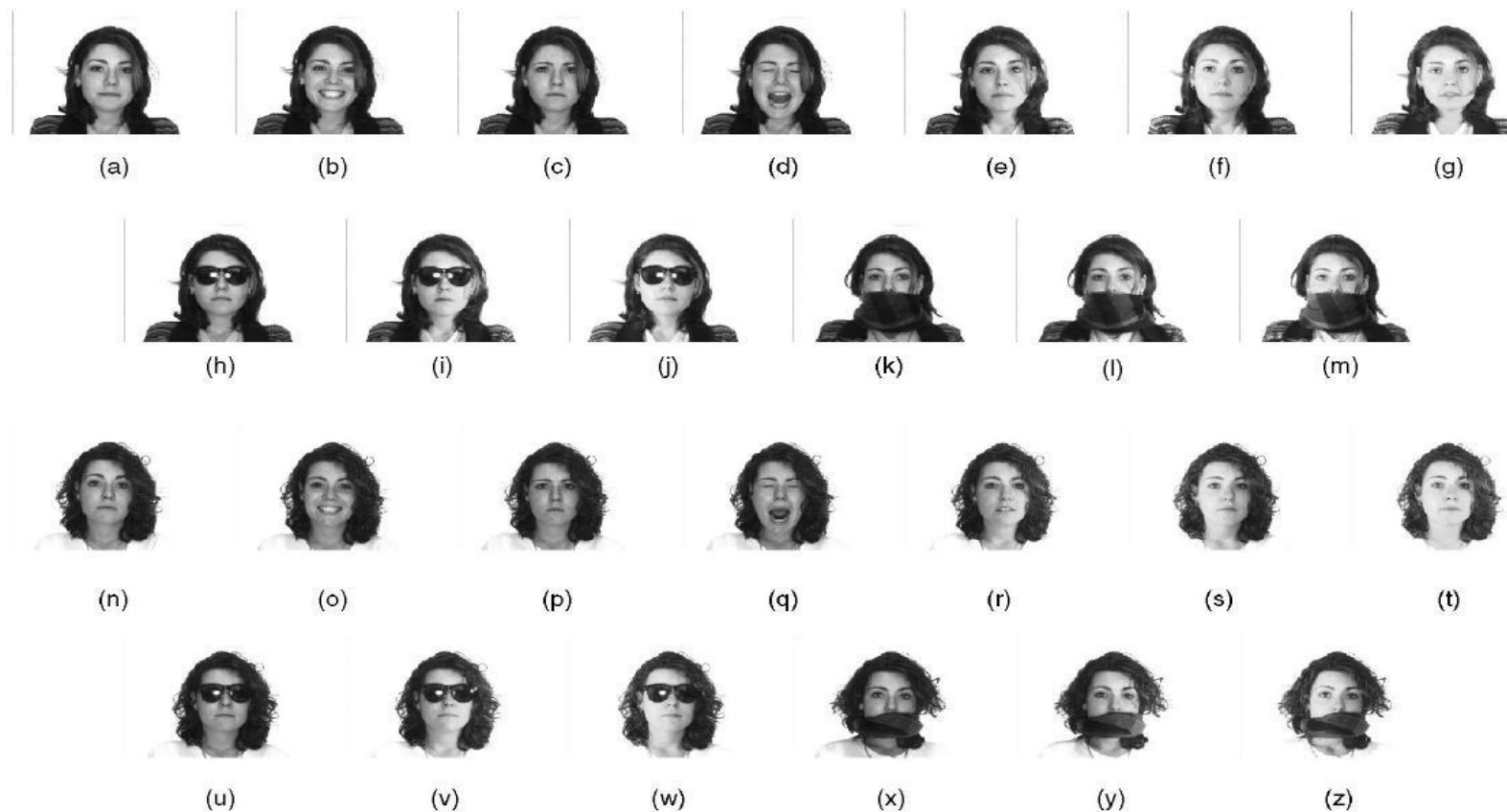
$$\sum n_i (\mu_i - \mu) = 0$$

Question

- Shall LDA always be **better** than PCA for classification task?

Is LDA always better than PCA?

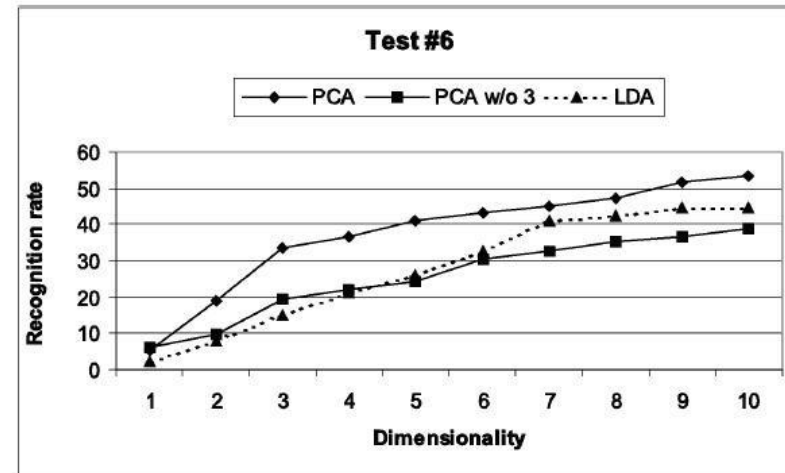
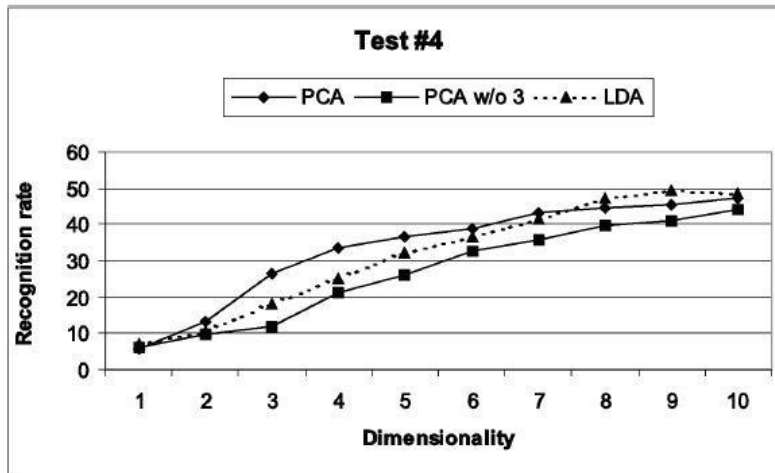
- Case Study: PCA versus LDA
 - A. Martinez, A. Kak, "PCA versus LDA", *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 23, no. 2, pp. 228-233, 2001.



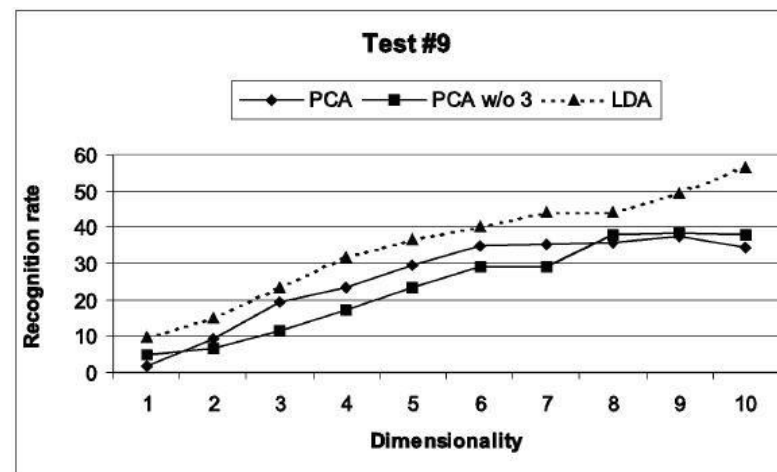
One subject in AR database

Is LDA always better than PCA?

- Small training dataset
(2 images/person for training and 5 images/person for testing)



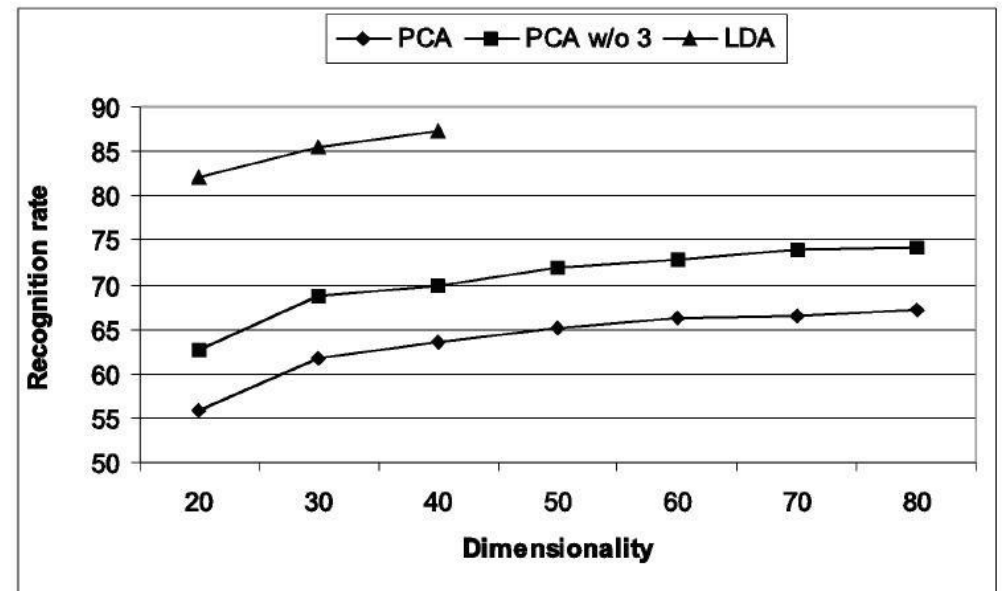
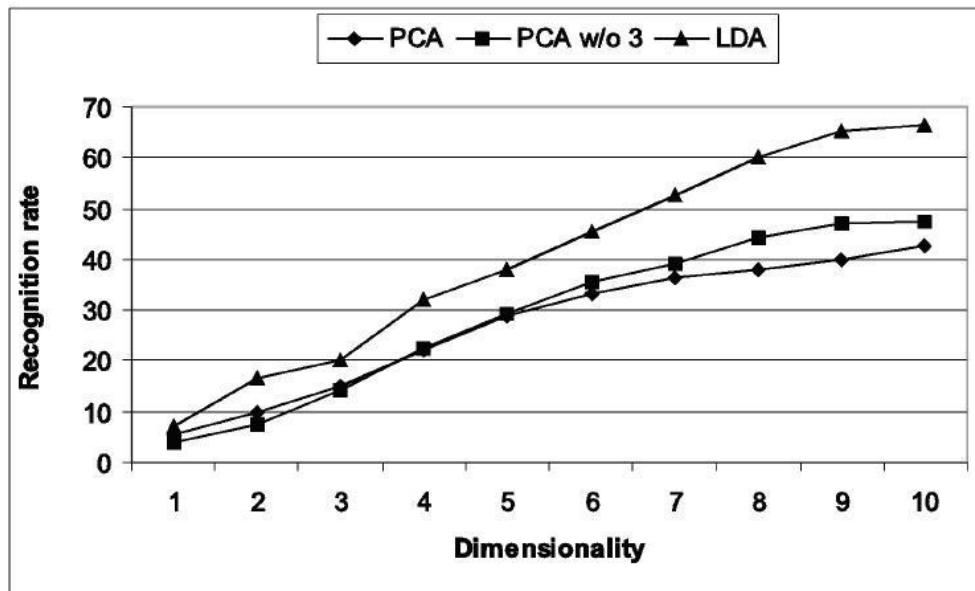
LDA is not always
better when the
sample size is small



PCA w/o 3: PCA without using the top 3 eigenvectors. #i index train/test splits

Is LDA always better than PCA?

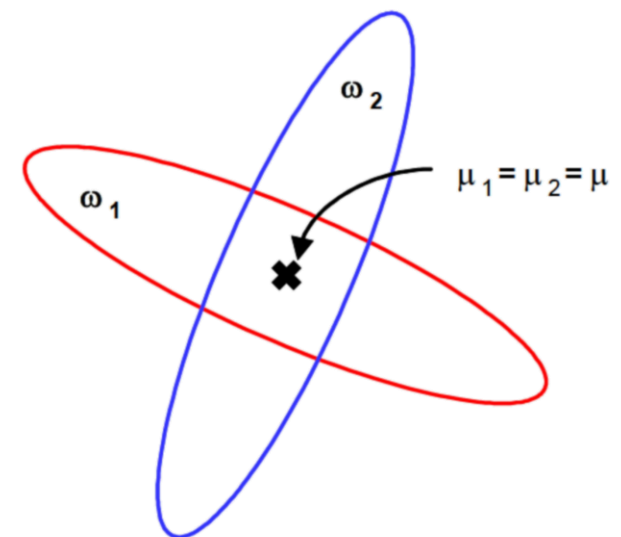
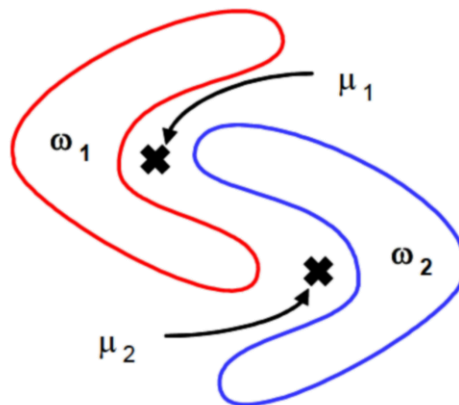
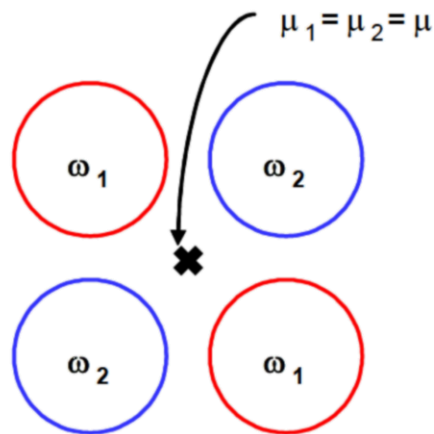
- Large training dataset
(13 images/person for training and 13 images/person for testing)



LDA outperforms PCA when the sample size is large

Limitations

- LDA produces at most $\min(d, C - 1)$ feature projections
 - If more features are needed, some other method must be employed to provide those additional features
- LDA is a parametric method (it assumes unimodal Gaussian likelihoods)
 - If the distributions are non-Gaussian, LDA projections may not preserve complex structure in the data needed for classification



A general framework:
Graph Embedding (GE)

GE: Problem Formulation of Feature Learning

- For a general feature learning problem

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Goal: Find a mapping function $F: \mathbf{x} \rightarrow \mathbf{y}$, where $\mathbf{y}_i \in \mathbb{R}^{d'}$, $d' \ll d$
(i.e., $\mathbf{y} = F(\mathbf{x})$)

- From the perspective of **Graph Embedding**

- ✓ Undirected Weighted Graph: $G = \{\mathbf{X}, \mathbf{S}\}$

Vertex set

Similarity (edge) matrix $\mathbf{S} \in \mathbb{R}^{N \times N}$

Dissimilar: negative

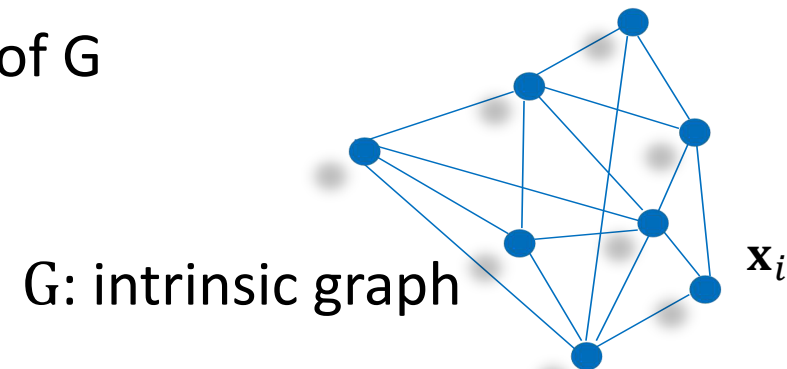
$\mathbf{S}_{i,i} = 0$

Similar: positive

- ✓ Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

Sum of similarities
between sample \mathbf{x}_i
and the other samples

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$



- ✓ Goal: Find the desired low-dimensional vector representations \mathbf{y} to **Preserve** these graph similarities.

GE: Problem Formulation of Feature Learning

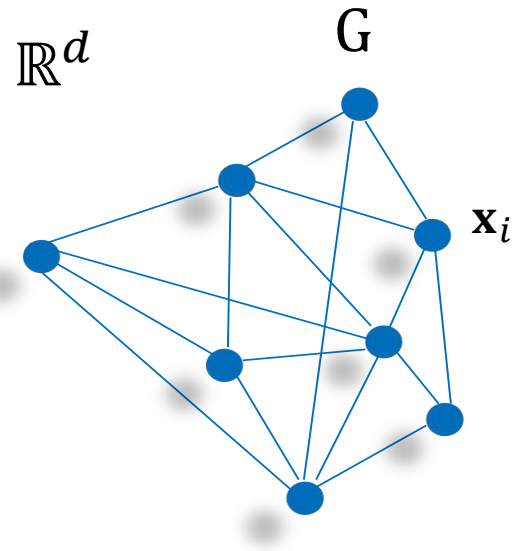
- Assume a simpler case (y is 1D)

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Intrinsic graph: $G = \{\mathbf{X}, \mathbf{S}\}$

Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$

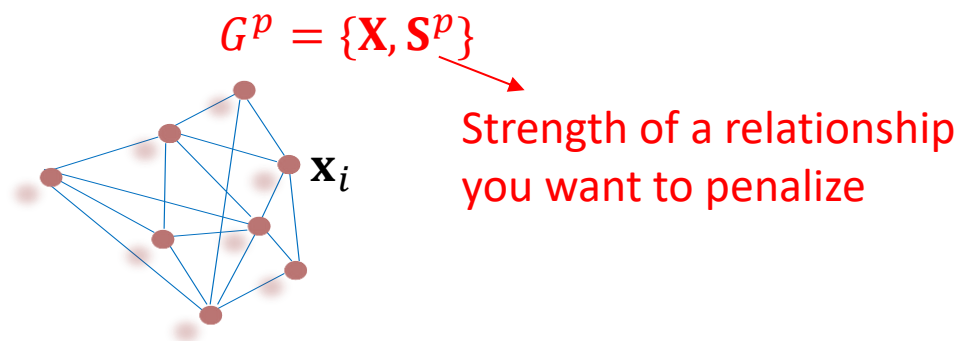


Goal: Find $\mathbf{y} \in \mathbb{R}$ to preserve graph similarity ($\mathbf{y} = [y_1, y_2, \dots, y_N]^T$)

$$\mathbf{y}^* = \operatorname{argmin}_{\mathbf{y}^T \mathbf{B} \mathbf{y} = c} \sum_{i \neq j} \|\mathbf{y}_i - \mathbf{y}_j\|^2 \mathbf{S}_{i,j}$$

c : constant,

\mathbf{B} : constraint matrix to avoid trivial solution



$$\mathbf{B} = \begin{cases} \text{A diagonal matrix: Scale normalization} \\ \text{Laplacian of a penalty graph } G^p: \\ \mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p \end{cases}$$

GE: Problem Formulation of Feature Learning

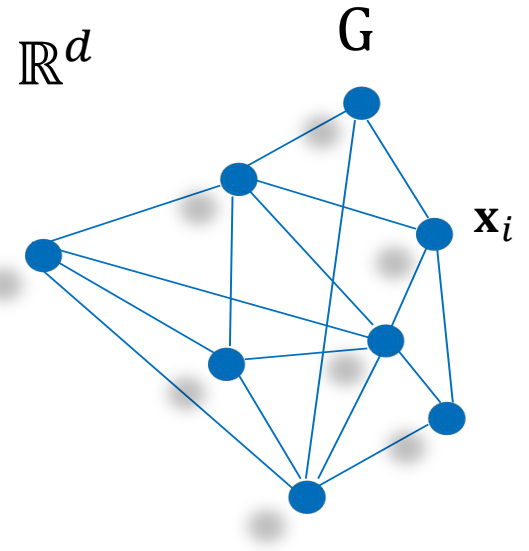
- Assume a simpler case (y is 1D)

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Intrinsic graph: $G = \{\mathbf{X}, \mathbf{S}\}$

The Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$



Goal: Find $y \in \mathbb{R}$ to preserve graph similarity ($y = [y_1, y_2, \dots, y_N]^T$)

$$y^* = \operatorname{argmin}_{y^T \mathbf{B} y = c} \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{S}_{i,j}$$

c : constant,

\mathbf{B} : constraint matrix to avoid trivial solution

$$\mathbf{B} = \begin{cases} \text{A diagonal matrix: Scale normalization} \\ \text{Laplacian of a penalty graph } G^p: \\ \mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p \end{cases}$$

$$y^T \mathbf{B} y = c \Leftrightarrow \sum_i \|y_i\|^2 \mathbf{B}_{i,i} = c$$

$$y^T \mathbf{B} y = c \Leftrightarrow \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{S}_{i,j}^p = 2c$$

Proof of Equivalence in Constraint

$$\mathbf{y}^T \mathbf{B} \mathbf{y} = c \Leftrightarrow \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{s}_{i,j}^p = 2c$$

$$\mathbf{D}_{i,i}^p = \sum_{i \neq j} \mathbf{s}_{i,j}^p$$



$$\begin{aligned} \mathbf{y}^T \mathbf{B} \mathbf{y} &= \mathbf{y}^T \mathbf{D}^p \mathbf{y} - \mathbf{y}^T \mathbf{S}^p \mathbf{y} \\ &= \sum_i \mathbf{D}_{i,i}^p y_i^2 - \sum_{i,j} y_i \mathbf{s}_{i,j}^p y_j \end{aligned}$$



$$\mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p$$

$$\begin{aligned} \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{s}_{i,j}^p &= \sum_{i \neq j} (y_i^2 - 2y_i y_j + y_j^2) \mathbf{s}_{i,j}^p \\ &= \sum_{i \neq j} y_i^2 \mathbf{s}_{i,j}^p - \sum_{i \neq j} 2y_i y_j \mathbf{s}_{i,j}^p + \sum_{i \neq j} y_j^2 \mathbf{s}_{i,j}^p \\ &= 2 \sum_i y_i^2 \sum_{j \neq i} \mathbf{s}_{i,j}^p - \sum_{i \neq j} 2y_i y_j \mathbf{s}_{i,j}^p \\ &= 2 \left(\sum_i y_i^2 \mathbf{D}_{i,i}^p - \sum_{i \neq j} y_i y_j \mathbf{s}_{i,j}^p \right) \end{aligned}$$

Since typically, no penalty should be posed when $i=j$ (i.e., $\mathbf{s}_{i,i}^p = 0$)

$$= 2 \left(\sum_i y_i^2 \mathbf{D}_{i,i}^p - \sum_{i,j} y_i y_j \mathbf{s}_{i,j}^p \right)$$

GE: Problem Formulation of Feature Learning

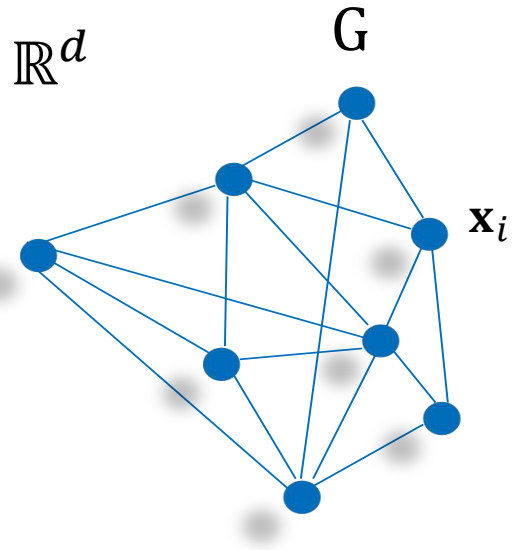
- Assume a simpler case (y is 1D)

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Intrinsic graph: $G = \{\mathbf{X}, \mathbf{S}\}$

The Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$



Goal: Find $\mathbf{y} \in \mathbb{R}$ to preserve graph similarity ($\mathbf{y} = [y_1, y_2, \dots, y_N]^T$)

$$\mathbf{y}^* = \operatorname{argmin}_{\mathbf{y}^T \mathbf{B} \mathbf{y} = c} \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{S}_{i,j}$$

c : constant,
 \mathbf{B} : constraint matrix to avoid trivial solution

Why this preserves graph similarity?

GE: Problem Formulation of Feature Learning

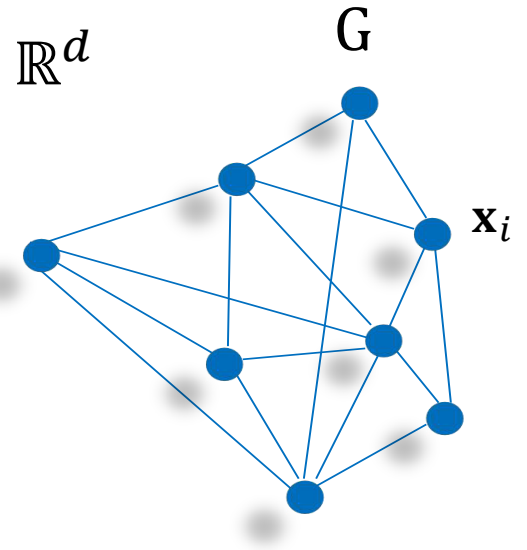
- Assume a simpler case (y is 1D)

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Intrinsic graph: $G = \{\mathbf{X}, \mathbf{S}\}$

The Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$



Goal: Find $\mathbf{y} \in \mathbb{R}$ to preserve graph similarity ($\mathbf{y} = [y_1, y_2, \dots, y_N]^T$)

$$\mathbf{y}^* = \underset{\mathbf{y}^T \mathbf{B} \mathbf{y} = c}{\operatorname{argmin}} \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{S}_{i,j} = \underset{\mathbf{y}^T \mathbf{B} \mathbf{y} = c}{\operatorname{argmin}} \mathbf{y}^T \mathbf{L} \mathbf{y}$$

c : constant,

\mathbf{B} : constraint matrix to avoid trivial solution

GE: Linearization

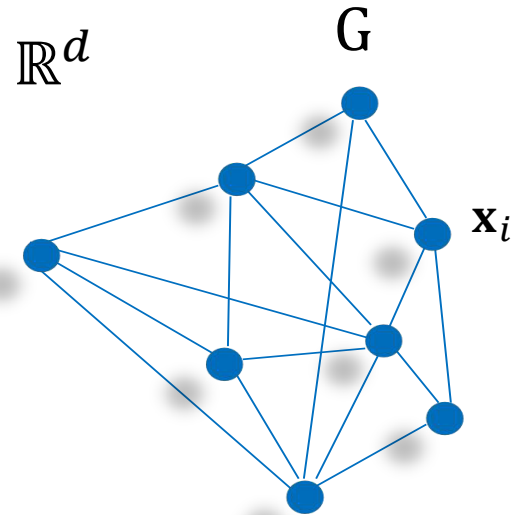
- Assume $\mathbf{y} = \mathbf{X}^T \mathbf{w}$

Data matrix for training: $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N]$, where $\mathbf{x}_i \in \mathbb{R}^d$

Intrinsic graph: $G = \{\mathbf{X}, \mathbf{S}\}$

The Diagonal matrix \mathbf{D} and the Laplacian matrix \mathbf{L} of G

$$\mathbf{D}_{i,i} = \sum_{i \neq j} \mathbf{S}_{i,j} \quad \mathbf{L} = \mathbf{D} - \mathbf{S}$$



Goal: Find $\mathbf{y} \in \mathbb{R}$ to preserve graph similarity ($\mathbf{y} = [y_1, y_2, \dots, y_N]^T$)

$$\mathbf{y}^* = \operatorname{argmin}_{\mathbf{y}^T \mathbf{B} \mathbf{y} = c} \sum_{i \neq j} \|y_i - y_j\|^2 \mathbf{S}_{i,j} = \operatorname{argmin}_{\mathbf{y}^T \mathbf{B} \mathbf{y} = c} \mathbf{y}^T \mathbf{L} \mathbf{y}$$



$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c} \sum_{i \neq j} \|\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j\|^2 \mathbf{S}_{i,j} = \operatorname{argmin}_{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}$$

GE: Linearization

- Assume $\mathbf{y} = \mathbf{X}^T \mathbf{w}$

$$\mathbf{w}^* = \underset{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c}{\operatorname{argmin}} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}$$

Lagrange Multipliers:

Given f to be optimized, subject to the constraint $g=0$, then, solve: $\nabla f = \lambda \nabla g$

$$f = \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}, \quad g = \mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} - c$$

Let λ be a Lagrange multiplier,

$$\nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}) = \lambda \nabla_{\mathbf{w}}(\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} - c)$$

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} = \lambda \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w}$$

Generalized eigenvalue decomposition problem:

$$\tilde{\mathbf{L}} \mathbf{w} = \lambda \tilde{\mathbf{B}} \mathbf{w}$$

Where $\tilde{\mathbf{L}} = \mathbf{X} \mathbf{L} \mathbf{X}^T$ and $\tilde{\mathbf{B}} = \mathbf{X} \mathbf{B} \mathbf{X}^T$

GE: Kernelization

- Extend to Nonlinear Cases

Map the data to a higher dimensional space: $\mathbf{x}_i \rightarrow \boldsymbol{\varphi}(\mathbf{x}_i)$

A kernel function K : $K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j)$

e.g., data \mathbf{x} is 2-D, i.e., $\mathbf{x} = (x_1, x_2)^T$

We map it to a 6-D space: $\boldsymbol{\varphi}(\mathbf{x}) = (x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1)^T$

$$\boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j) = (x_{i,1}^2, x_{i,2}^2, \sqrt{2}x_{i,1}x_{i,2}, \sqrt{2}x_{i,1}, \sqrt{2}x_{i,2}, 1) \begin{pmatrix} x_{j,1}^2 \\ x_{j,2}^2 \\ \sqrt{2}x_{j,1}x_{j,2} \\ \sqrt{2}x_{j,1} \\ \sqrt{2}x_{j,2} \\ 1 \end{pmatrix}$$

d -dimensional \rightarrow m -dimensional

No $K()$: $O(m)$

With $K()$: $O(d)$

Kernel saves computation!

$$= x_{i,1}^2 x_{j,1}^2 + x_{i,2}^2 x_{j,2}^2 + 2x_{i,1}x_{i,2}x_{j,1}x_{j,2} + 2x_{i,1}x_{j,1} + 2x_{i,2}x_{j,2} + 1$$

$$= (x_{i,1}x_{j,1} + x_{i,2}x_{j,2} + 1)^2$$

$$= (\mathbf{x}_i^T \mathbf{x}_j + 1)^2 = K(\mathbf{x}_i, \mathbf{x}_j)$$

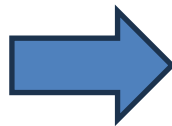
GE: Kernelization

- Extend to Nonlinear Cases

Map the data to a higher dimensional space: $\mathbf{x}_i \rightarrow \boldsymbol{\varphi}(\mathbf{x}_i)$

A kernel function K : $K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^T \boldsymbol{\varphi}(\mathbf{x}_j)$

$$\mathbf{w}^* = \underset{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c}{\operatorname{argmin}} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}$$



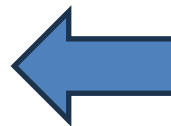
$$\mathbf{w}^* = \underset{\mathbf{w}^T \boldsymbol{\Phi} \mathbf{B} \boldsymbol{\Phi}^T \mathbf{w} = c}{\operatorname{argmin}} \mathbf{w}^T \boldsymbol{\Phi} \mathbf{L} \boldsymbol{\Phi}^T \mathbf{w}$$

where $\boldsymbol{\Phi} = [\boldsymbol{\varphi}(\mathbf{x}_1), \dots, \boldsymbol{\varphi}(\mathbf{x}_N)]$

Assume projection vector $\mathbf{w} = \sum_i \alpha_i \boldsymbol{\varphi}(\mathbf{x}_i) = \boldsymbol{\Phi} \boldsymbol{\alpha}$



$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{B} \boldsymbol{\Phi}^T \boldsymbol{\Phi} \boldsymbol{\alpha} = c}{\operatorname{argmin}} \boldsymbol{\alpha}^T \boldsymbol{\Phi}^T \boldsymbol{\Phi} \mathbf{L} \boldsymbol{\Phi}^T \boldsymbol{\Phi} \boldsymbol{\alpha}$$



$$\boldsymbol{\alpha}^* = \underset{\boldsymbol{\alpha}^T \mathbf{G} \mathbf{B} \mathbf{G} \boldsymbol{\alpha} = c}{\operatorname{argmin}} \boldsymbol{\alpha}^T \mathbf{G} \mathbf{L} \mathbf{G} \boldsymbol{\alpha}$$

$$\text{Since } \boldsymbol{\Phi}^T \boldsymbol{\Phi} = \begin{bmatrix} \boldsymbol{\varphi}(\mathbf{x}_1)^T \boldsymbol{\varphi}(\mathbf{x}_1) & \cdots & \boldsymbol{\varphi}(\mathbf{x}_1)^T \boldsymbol{\varphi}(\mathbf{x}_N) \\ \vdots & \ddots & \vdots \\ \boldsymbol{\varphi}(\mathbf{x}_N)^T \boldsymbol{\varphi}(\mathbf{x}_1) & \cdots & \boldsymbol{\varphi}(\mathbf{x}_N)^T \boldsymbol{\varphi}(\mathbf{x}_N) \end{bmatrix}$$

Kernel Gram matrix \mathbf{G}

What is the solution?

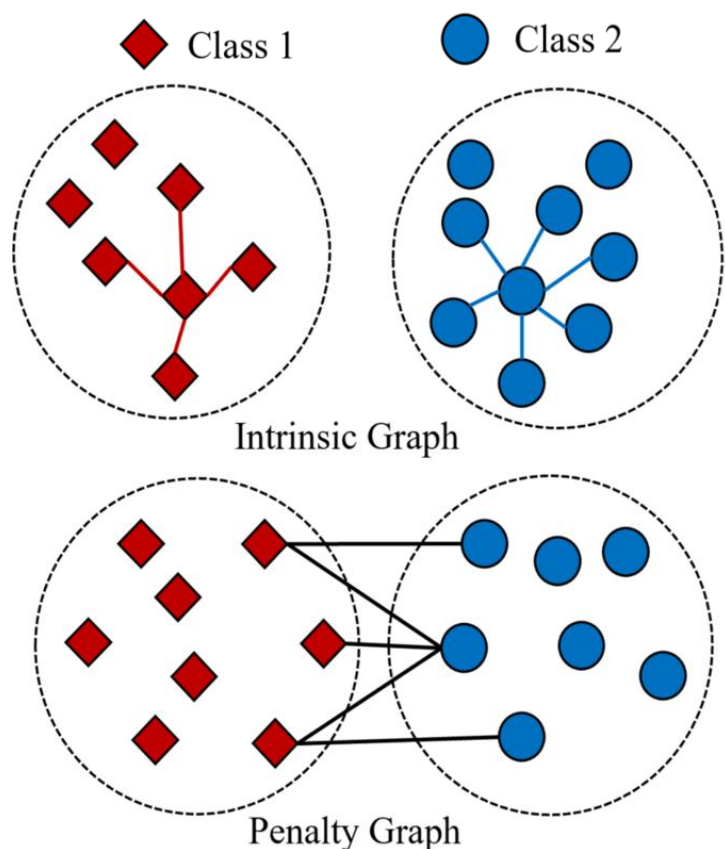
A New Algorithm using GE: Marginal Fisher Analysis (MFA)

New Algorithm: Marginal Fisher Analysis

- **Motivation**

-Avoid certain limitations of LDA in terms of the **data distribution assumption** and **available projection direction**

- **Utilize GE**



- ✓ **Intrinsic Graph G : Intraclass compactness**
(Each sample is connected to its k_1 -nearest neighbors of the same class)

$S_{i,j} = S_{j,i} = 1$: if for each sample \mathbf{x}_i , \mathbf{x}_j is among the k_1 -nearest neighbors of \mathbf{x}_i in the same class

- ✓ **Penalty Graph G^p : Interclass separability**
(Marginal point pairs of different classes are connected)

$S_{i,j}^p = 1$: if for each class c_l , the pair $(\mathbf{x}_i, \mathbf{x}_j)$ is among the k_2 -shortest pairs in the set $\{(\mathbf{x}_i, \mathbf{x}_j), \mathbf{x}_i \in \mathbf{X}_{cl}, \mathbf{x}_j \notin \mathbf{X}_{cl}\}$

New Algorithm: Marginal Fisher Analysis

- Intraclass Compactness following Linearization of GE**

$$\begin{aligned}\widehat{S}_c &= \sum_{i \neq j} \|\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j\|^2 \mathbf{S}_{i,j} = 2\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} \\ &= \sum_i \sum_{i \in N_{k_1}^+(j) \text{ or } j \in N_{k_1}^+(i)} \|\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j\|^2 = 2\mathbf{w}^T \mathbf{X} (\mathbf{D} - \mathbf{S}) \mathbf{X}^T \mathbf{w}\end{aligned}$$

$N_{k_1}^+(i)$ is the index set of the k_1 -nearest neighbors of the sample \mathbf{x}_i in the same class

$$\mathbf{S}_{i,j} = \begin{cases} 1, & \text{If } i \in N_{k_1}^+(j) \text{ or } j \in N_{k_1}^+(i) \\ 0, & \text{else} \end{cases}$$

- Interclass Separability following Linearization of GE**

$$\begin{aligned}\widehat{S}_p &= \sum_{i \neq j} \|\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j\|^2 \mathbf{S}_{i,j}^p = 2\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w} \\ &= \sum_i \sum_{(i,j) \in P_{k_2}(cl_i) \text{ or } (i,j) \in P_{k_2}(cl_j)} \|\mathbf{w}^T \mathbf{x}_i - \mathbf{w}^T \mathbf{x}_j\|^2 = 2\mathbf{w}^T \mathbf{X} (\mathbf{D}^p - \mathbf{S}^p) \mathbf{X}^T \mathbf{w}\end{aligned}$$

$P_{k_2}(cl)$ is a set of data pairs that are the k_2 nearest pairs among the set $\{(i,j), i \in \text{class } cl, j \notin \text{class } cl\}$

$$\mathbf{S}_{i,j}^p = \begin{cases} 1, & \text{If } (i,j) \in P_{k_2}(cl_i) \text{ or } (i,j) \in P_{k_2}(cl_j) \\ 0, & \text{else} \end{cases}$$

New Algorithm: Marginal Fisher Analysis

- Marginal Fisher Criterion

Minimize Intraclass Compactness \hat{S}_c , Maximize Interclass Separability \hat{S}_p

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}(\mathbf{D} - \mathbf{S})\mathbf{X}^T \mathbf{w}}{\mathbf{w}^T \mathbf{X}(\mathbf{D}^p - \mathbf{S}^p)\mathbf{X}^T \mathbf{w}}$$

Which is a special linearization of GE with $\mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p$



New Algorithm: Marginal Fisher Analysis

General Linearization:

MFA:

$$\mathbf{w}^* = \underset{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c}{\operatorname{argmin}} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} \quad \begin{array}{c} \text{=====} \\ \text{=====} \end{array} \quad \mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}}{\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w}}$$

Solve $\tilde{\mathbf{L}} \mathbf{w} = \lambda \tilde{\mathbf{B}} \mathbf{w}$,

where $\tilde{\mathbf{L}} = \mathbf{X} \mathbf{L} \mathbf{X}^T$ and $\tilde{\mathbf{B}} = \mathbf{X} \mathbf{B} \mathbf{X}^T$

$$\frac{d}{d\mathbf{w}} J(\mathbf{w}) = 0$$

$$\frac{(\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w}) \frac{d}{d\mathbf{w}} (\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}) - (\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}) \frac{d}{d\mathbf{w}} (\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w})}{(\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w})^2} = 0$$

$$2(\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w}) \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} - 2(\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}) \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w} = 0$$

Divide by

$\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w}$

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} - \frac{(\mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w})}{(\mathbf{w}^T \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w})} \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w} = 0$$

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} - J(\mathbf{w}) \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w} = 0$$

$$\mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w} = J(\mathbf{w}) \mathbf{X} \mathbf{L}^p \mathbf{X}^T \mathbf{w}$$

New Algorithm: Marginal Fisher Analysis

- Marginal Fisher Criterion

Minimize Intraclass Compactness \widehat{S}_c , Maximize Interclass Separability \widehat{S}_p

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}} \frac{\mathbf{w}^T \mathbf{X}(\mathbf{D} - \mathbf{S})\mathbf{X}^T \mathbf{w}}{\mathbf{w}^T \mathbf{X}(\mathbf{D}^p - \mathbf{S}^p)\mathbf{X}^T \mathbf{w}}$$

Which is a special linearization of GE with $\mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p$

General Linearization:

$$\mathbf{w}^* = \operatorname{argmin}_{\mathbf{w}^T \mathbf{X} \mathbf{B} \mathbf{X}^T \mathbf{w} = c} \mathbf{w}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{w}$$

If $\mathbf{B} = \mathbf{L}^p = \mathbf{D}^p - \mathbf{S}^p$, the same as marginal fisher criterion!

Experiments: Face Recognition



Sample Images from CMU PIE database (after cropping)

	G4/P16	G5/P15	G6/P14
PCA	43.4% (252)	49.2% (264)	49.2% (341)
Fisherface	76.9% (62)	83.6% (61)	88.9% (62)
MFA	82.5% (70)	87.3% (117)	90.5% (69)
PCA+LDA	78.8% (94, 44)	84.3% (93, 36)	91.2% (91, 20)
PCA+MFA	83.9% (99, 71)	87.5% (99, 104)	90.8% (91, 32)

Gm/Pn : m images per person are randomly selected for training and the remaining n images are used for testing.

Papers to Read and Study

- Shuicheng Yan, Dong Xu et al.: Graph Embedding and Extensions: A General Framework for Dimensionality Reduction. T-PAMI 2007.
- Jieping Ye, Ravi Janardan, Cheonghee Park, and Haesun Park. An optimization criterion for generalized discriminant analysis on undersampled problems. IEEE Transactions on Pattern Analysis and Machine Intelligence. Vol. 26, No. 8, pp. 982—994, 2004.
- A. Martinez, A. Kak, "PCA versus LDA", IEEE Transactions on Pattern Analysis and Machine Intelligence, vol. 23, no. 2, pp. 228-233, 2001.