

Machine Learning & Pattern Recognition

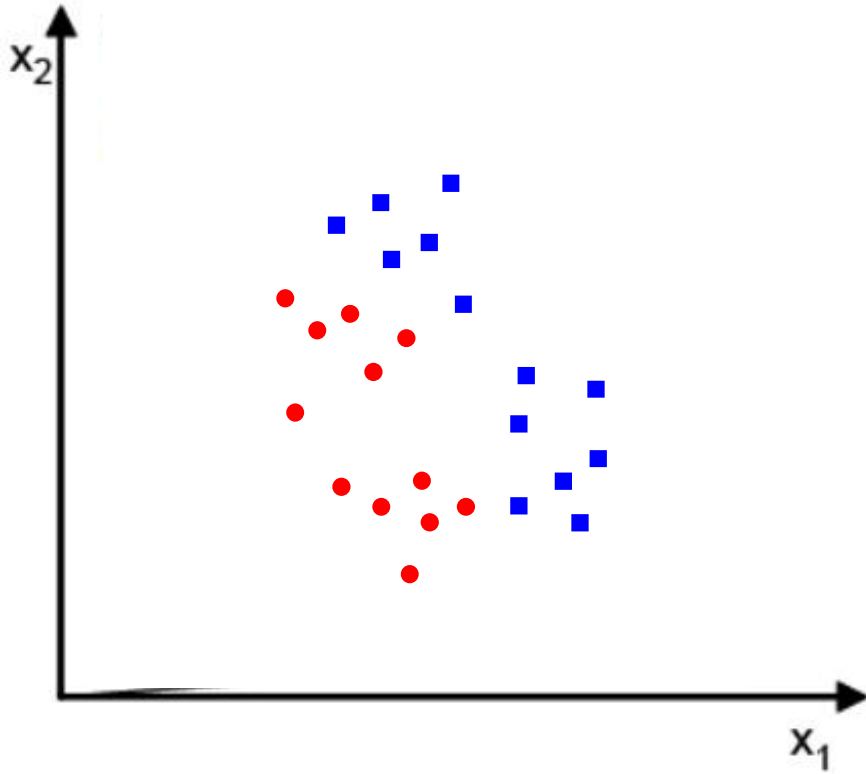
SONG Xuemeng

sxmustc@gmail.com

<http://xuemeng.bitcron.com/>

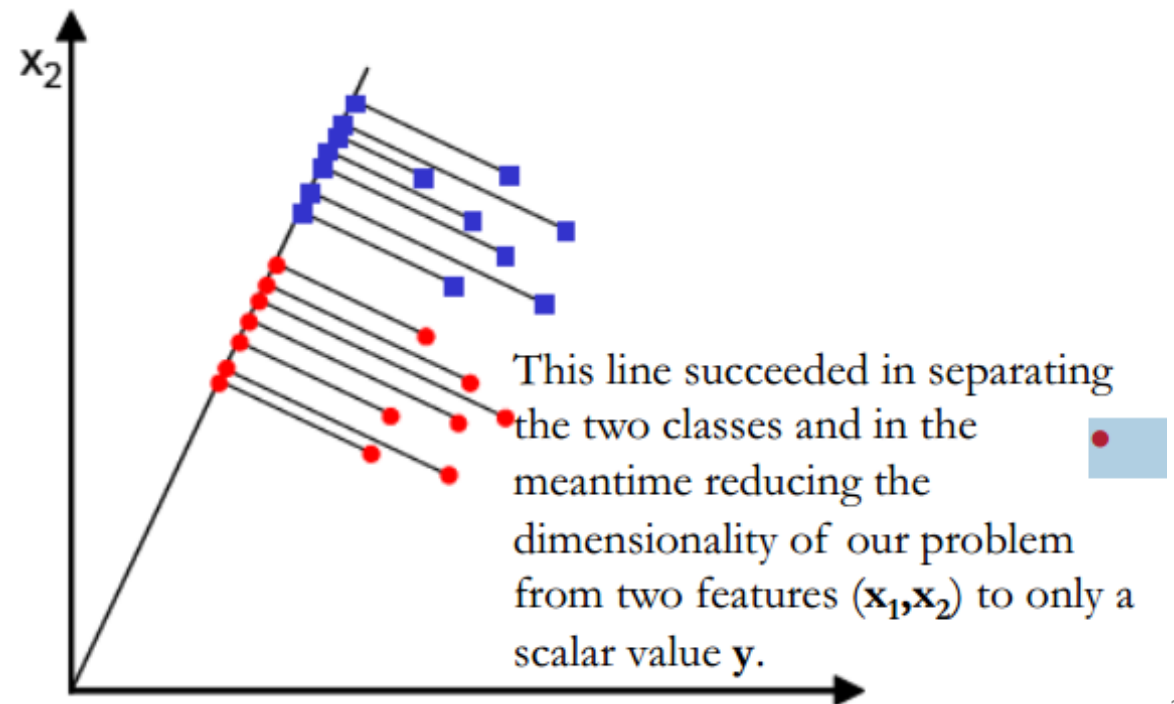
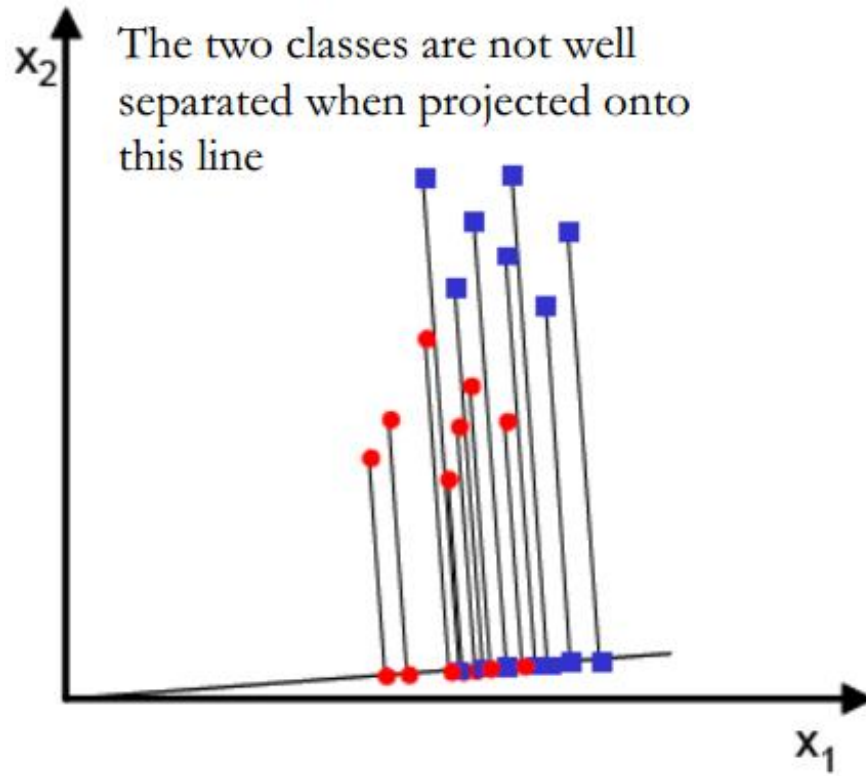
What is a Good Projection?

- Given a set of points (2-d) from two classes, we want to project them to a line that can well separate them.
- What is a good criterion?



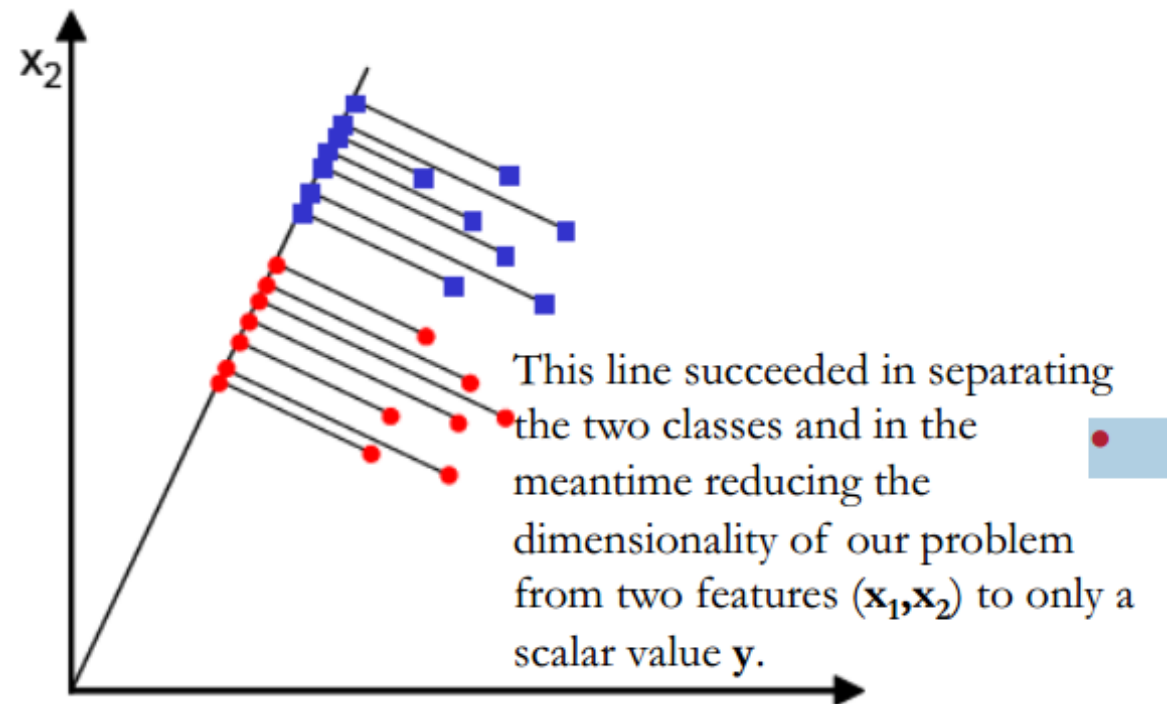
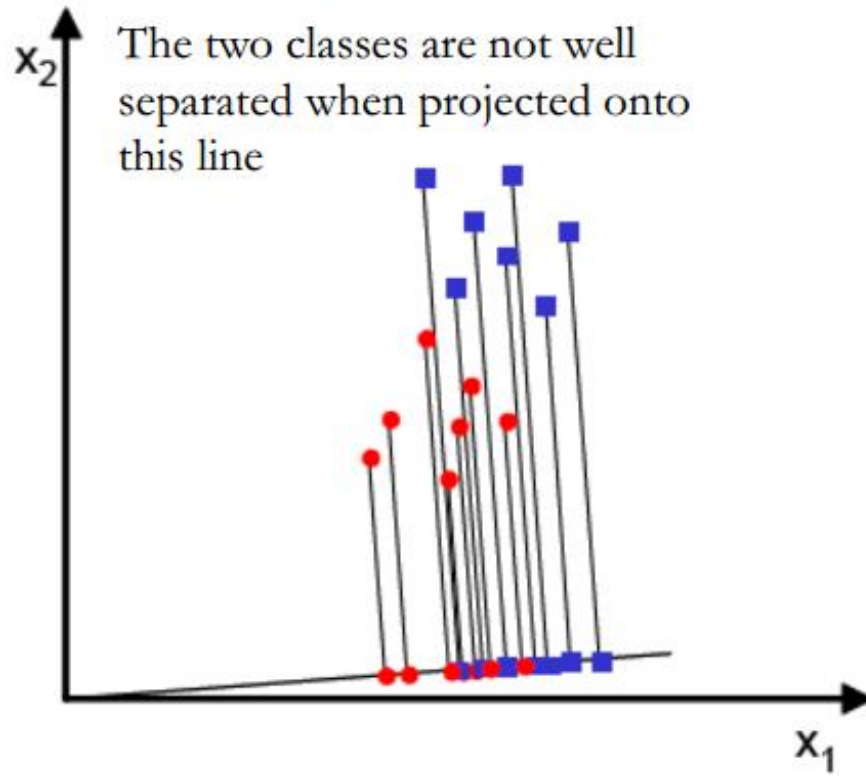
What is a Good Projection?

- What is a good criterion?
 - Separating different classes



What is a Good Projection?

- What is a good criterion?
 - Separating different classes
 - Maximize the between-class distance (means)



What is a Good Projection?

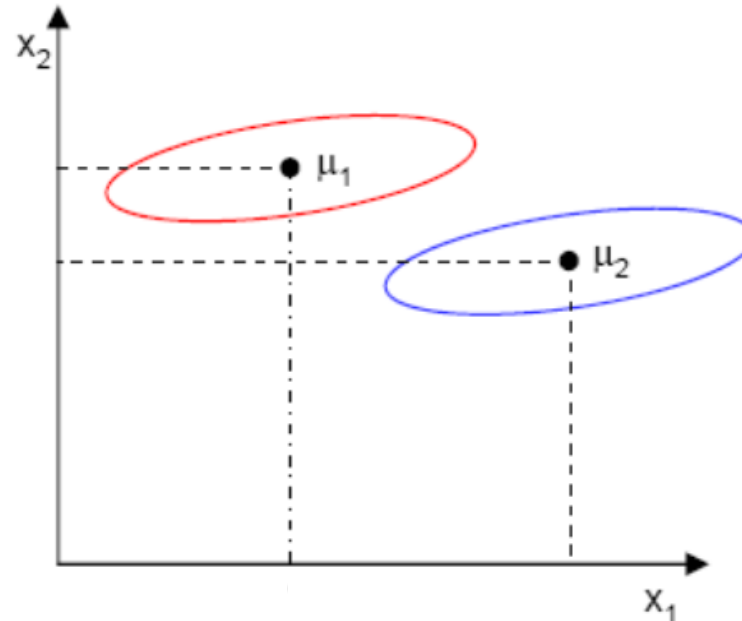
- What is a good criterion?
 - Separating different classes
 - Maximize the between-class distance (means)

Is it enough?

What is a Good Projection?

- What is a good criterion?
 - Separating different classes
 - Maximize the between-class distance (means)

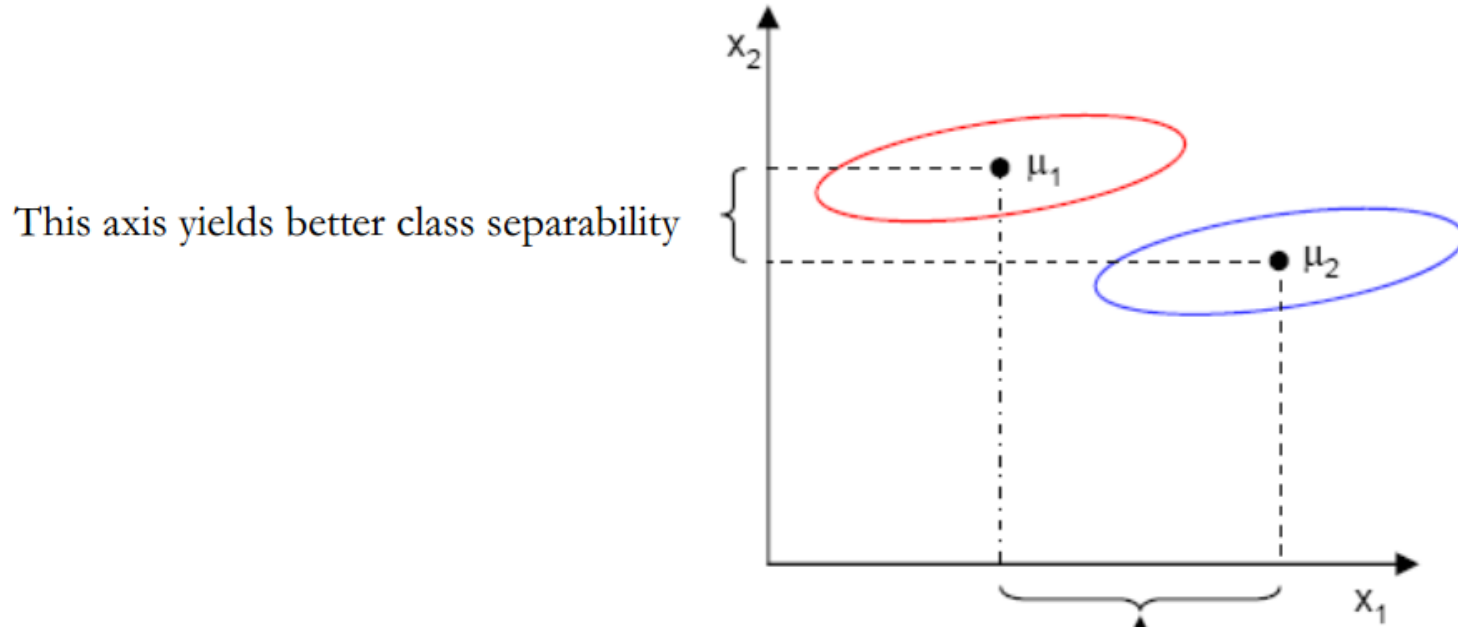
Is it enough?



What is a Good Projection?

- What is a good criterion?
 - Separating different classes
 - Maximize the between-class distance (means)

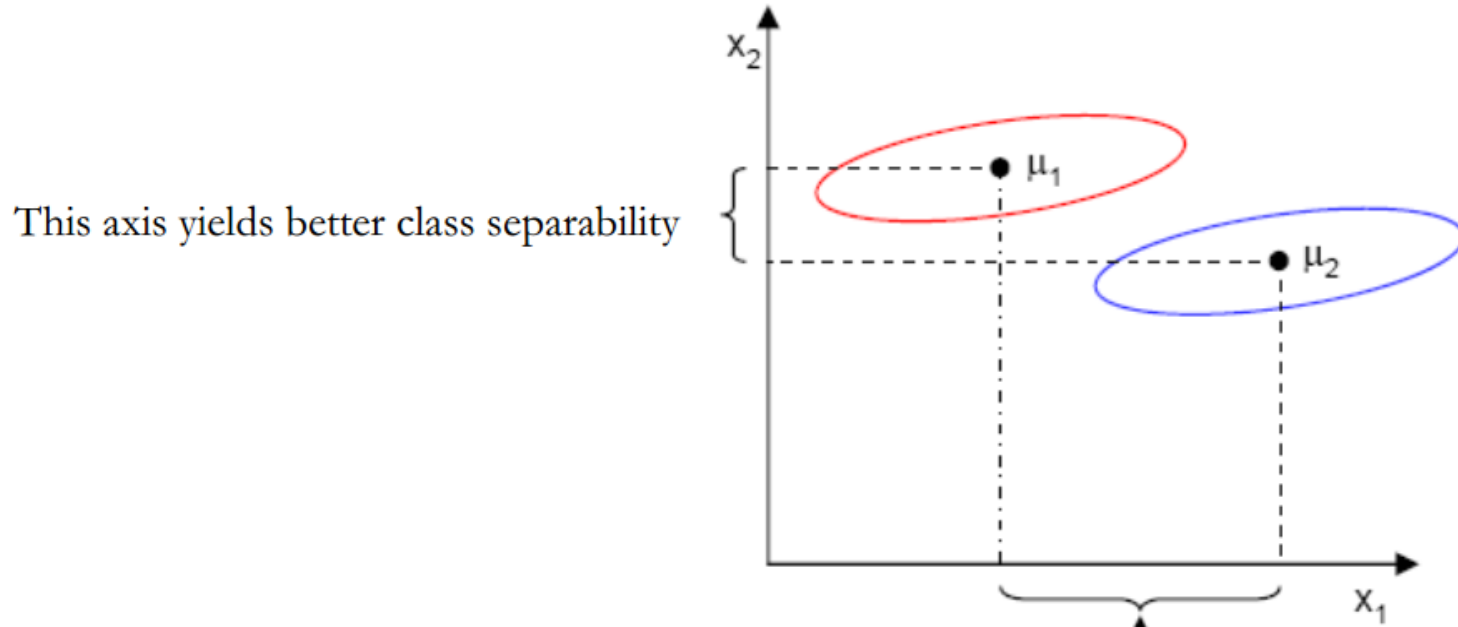
Is it enough?



This axis has a larger distance between means

What is a Good Projection?

- What is a good criterion?
 - Separating different classes
 - Maximize the between-class distance (means)
 - Minimize the within-class variability (scatter)



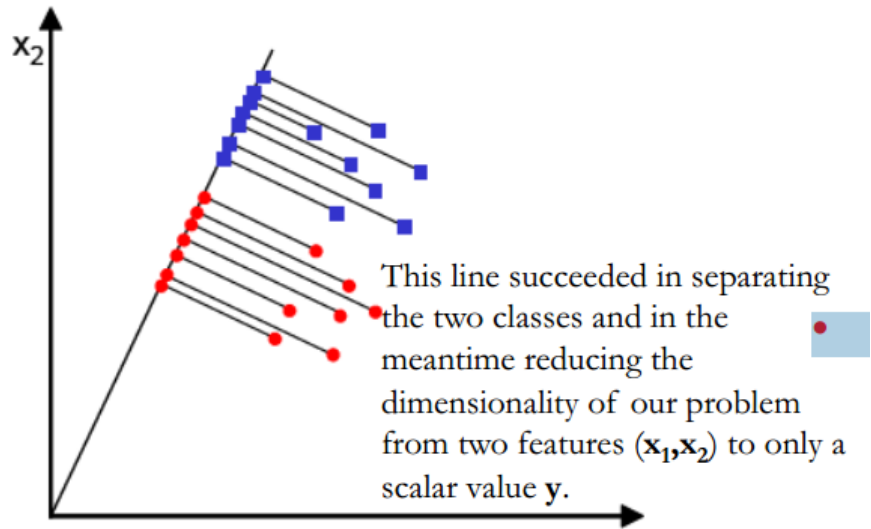
This axis has a larger distance between means

Linear Discriminant Analysis

- We have N d -dimensional samples from C classes, e.g., seabass, tuna, ...
- Each class has n_i samples, where $i = 1, 2, \dots, C$
- Stacking these samples from different classes into one big fat matrix $\mathbf{X} \in \mathbb{R}^{d \times N}$ such that each column represents one sample $\mathbf{x} \in \mathbb{R}^{d \times 1}$.
- We seek to obtain a transformation to project the d -dimensional samples in \mathbf{X} onto a p -dimensional subspace ($p < d$), such that after the projection we have:

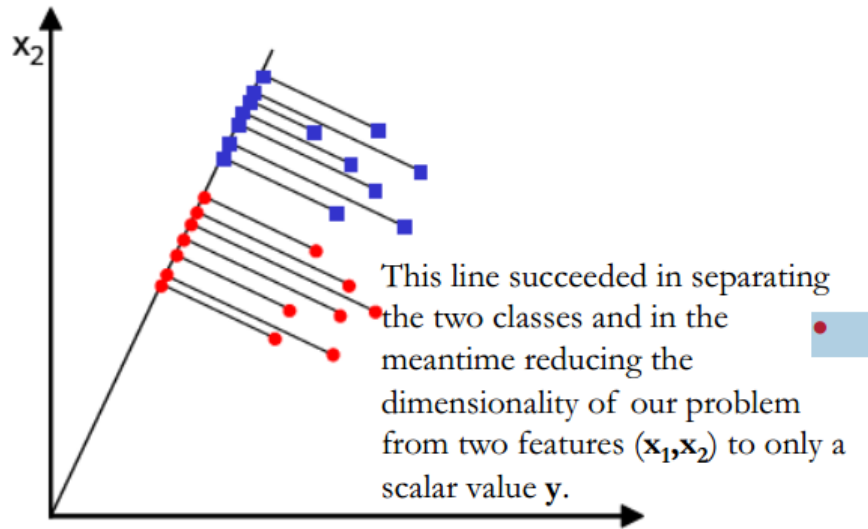
class means to be as far apart from each other as possible	→	the between-class scatter to be large
samples from the same class to be as close to their mean as possible	→	the within-class scatter to be small

Linear Discriminant Analysis



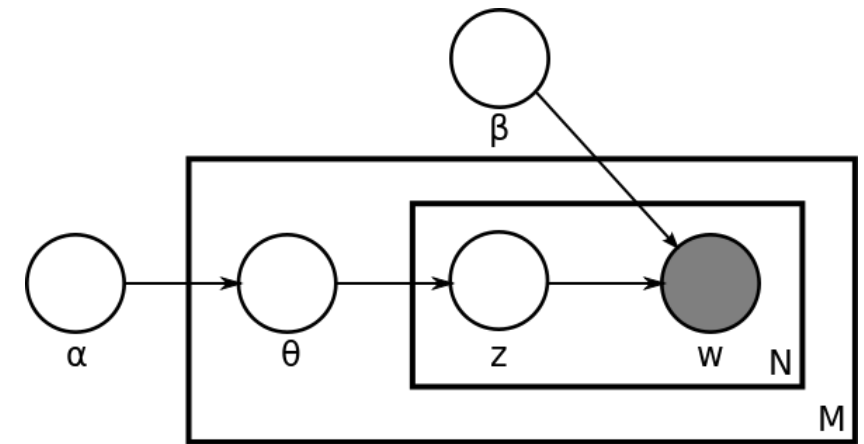
Linear Discriminant Analysis, a method to find a linear combination of features that **separates** two or more classes of objects.

Linear Discriminant Analysis



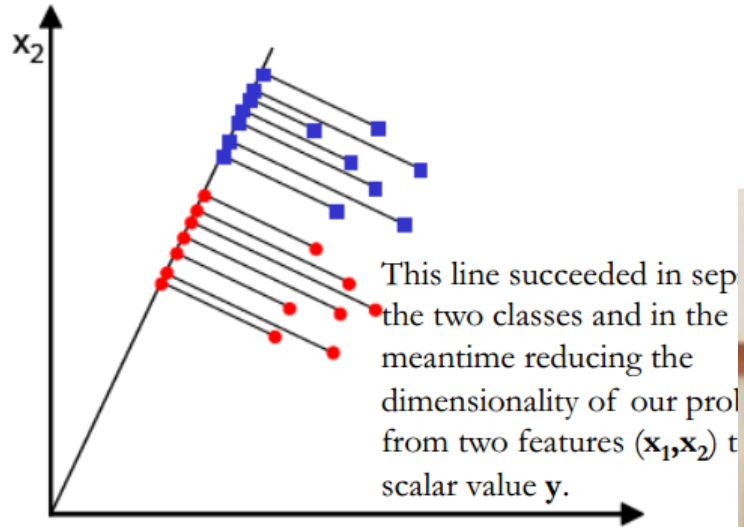
Linear Discriminant Analysis, a method to find a linear combination of features that **separates** two or more classes of objects.

Latent Dirichlet Allocation



In natural language processing, latent Dirichlet allocation (LDA) is an example of a **topic** model.
https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation

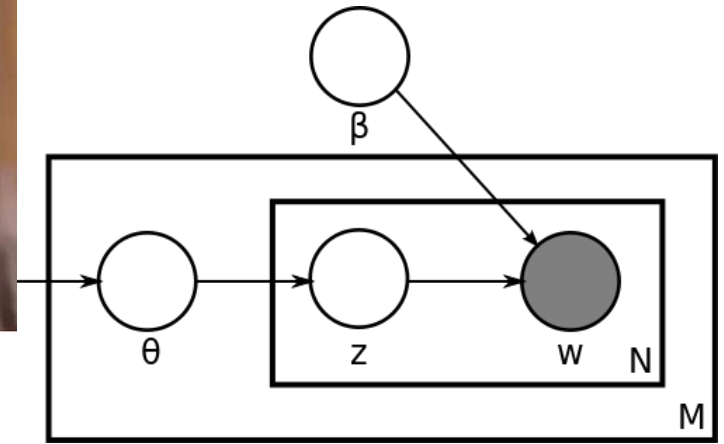
Linear Discriminant Analysis



Linear Discriminant Analysis, or LDA, is a technique to find a linear combination of features that **separates** two or more classes of objects.



Latent Dirichlet Allocation



In natural language processing, latent Dirichlet allocation (LDA) is an example of a **topic** model.
https://en.wikipedia.org/wiki/Latent_Dirichlet_allocation

Linear Discriminant Analysis

- Linear Discriminant Analysis—Two Classes
- Linear Discriminant Analysis—C Classes

Linear Discriminant Analysis—Two Classes

- Assume we have d -dimensional samples $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, n_1 of which belong to C_1 and n_2 belong to C_2 .
- We seek to obtain a transformation $\boldsymbol{\theta} \in \mathbb{R}^{d \times 1}$ that projects the samples \mathbf{x} onto a line ($p = 1$).

- $y_i = \boldsymbol{\theta}^T \mathbf{x}_i$, where $\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ \vdots \\ x_{id} \end{bmatrix}$ and $\boldsymbol{\theta} = \begin{bmatrix} \theta_1 \\ \vdots \\ \theta_d \end{bmatrix}$
- where $\boldsymbol{\theta}$ is the projection vectors used to project \mathbf{x} to y .

Statistical Facts

Within-class scatter:

$$\mathbf{S}_w = \sum_{x \in C_1} (\mathbf{x} - \boldsymbol{\mu}_1)(\mathbf{x} - \boldsymbol{\mu}_1)^T + \sum_{x \in C_2} (\mathbf{x} - \boldsymbol{\mu}_2)(\mathbf{x} - \boldsymbol{\mu}_2)^T \quad \mathbf{S}_w \in \mathbb{R}^{d \times d}$$

Between-class scatter:

$$\mathbf{S}_b = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \quad \mathbf{S}_b \in \mathbb{R}^{d \times d}$$

Class mean vector (sample):

$$\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{x \in C_i} \mathbf{x}, \boldsymbol{\mu}_i \in \mathbb{R}^{d \times 1}$$

Linear Discriminant Analysis—Two Classes

- The mean vector of each class in \mathbf{x} and y feature space is:

$$\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in C_i} \mathbf{x} \qquad \tilde{\boldsymbol{\mu}}_i = \frac{1}{n_i} \sum_{y \in C_i} y = \frac{1}{n_i} \sum_{\mathbf{x} \in C_i} \boldsymbol{\theta}^T \mathbf{x} = \boldsymbol{\theta}^T \boldsymbol{\mu}_i$$

- Projecting \mathbf{x} to y will lead to projecting the mean of \mathbf{x} to the mean of y .

- The within-class scatter: $\tilde{\mathbf{S}}_w = \sum_{y \in C_1} (y - \tilde{\mu}_1)^2 + \sum_{y \in C_2} (y - \tilde{\mu}_2)^2 = \boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta}$

- The between-class scatter: $\tilde{\mathbf{S}}_b = (\tilde{\mu}_1 - \tilde{\mu}_2)^2 = \boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta}$

Linear Discriminant Analysis—Two Classes

- On one hand, we want maximize the distance between the projected means:


$$\begin{aligned} J_1(\boldsymbol{\theta}) &= (\tilde{\mu}_1 - \tilde{\mu}_2)^2 = (\boldsymbol{\theta}^T \boldsymbol{\mu}_1 - \boldsymbol{\theta}^T \boldsymbol{\mu}_2)^2 \\ &= \boldsymbol{\theta}^T (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\theta} = \boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} = \tilde{\mathbf{S}}_b \end{aligned}$$

- On the other hand, we want minimize the within-class scatter:

$$J_2(\boldsymbol{\theta}) = \tilde{\mathbf{S}}_{w1} + \tilde{\mathbf{S}}_{w2} = \tilde{\mathbf{S}}_w = \boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta}$$

Linear Discriminant Analysis—Two Classes

- We can finally express the Fisher criterion in terms of \mathbf{S}_w and \mathbf{S}_b :


$$\begin{aligned}\max_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) &= \frac{\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta}}{\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta}} \\ \min_{\boldsymbol{\theta}} & -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} \\ \text{s.t. } & \boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} = 1\end{aligned}$$

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = ?$$

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

- $\boldsymbol{\theta}$: the eigenvectors of $\mathbf{S}_w^{-1} \mathbf{S}_b$, and λ is the corresponding eigenvalue.
- How to choose $\boldsymbol{\theta}$?

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

- Remember the objective function

$$\begin{cases} \min_{\boldsymbol{\theta}} -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} \\ \text{s.t. } \boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} = 1 \end{cases}$$

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

- Remember the objective function

$$\begin{cases} \min_{\boldsymbol{\theta}} -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} \\ \text{s.t. } \boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} = 1 \end{cases} \quad \begin{matrix} \mathbf{S}_b \boldsymbol{\theta}^* = \lambda \mathbf{S}_w \boldsymbol{\theta}^* \\ \Rightarrow \end{matrix} \quad -\boldsymbol{\theta}^{*T} \mathbf{S}_b \boldsymbol{\theta}^* = -\lambda \boldsymbol{\theta}^{*T} \mathbf{S}_w \boldsymbol{\theta}^* = -\lambda$$

- How to choose? The eigenvector corresponds to the **largest** eigenvalue.

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

- Alternatively, as $\mathbf{S}_b = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$,

$$\mathbf{S}_b \boldsymbol{\theta} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\theta}$$

Linear Discriminant Analysis—Two Classes

- Let λ be a **Lagrange multiplier**

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} + \lambda(\boldsymbol{\theta}^T \mathbf{S}_w \boldsymbol{\theta} - 1)$$

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = -2\mathbf{S}_b \boldsymbol{\theta} + 2\lambda \mathbf{S}_w \boldsymbol{\theta} = 0 \quad \Rightarrow \quad \mathbf{S}_b \boldsymbol{\theta} = \lambda \mathbf{S}_w \boldsymbol{\theta}$$

- As $\mathbf{S}_b = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$, $\mathbf{S}_b \boldsymbol{\theta} = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \boxed{(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\theta}}$
- Let $\mathbf{S}_b \boldsymbol{\theta} = \lambda_{\boldsymbol{\theta}}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ then $\lambda \mathbf{S}_w \boldsymbol{\theta} = \lambda_{\boldsymbol{\theta}}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ $\lambda_{\boldsymbol{\theta}}$
- The scale of $\boldsymbol{\theta}^*$ does not matter, only direction matters.

$$\boldsymbol{\theta}^* = \mathbf{S}_w^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Linear Discriminant Analysis—Two Classes

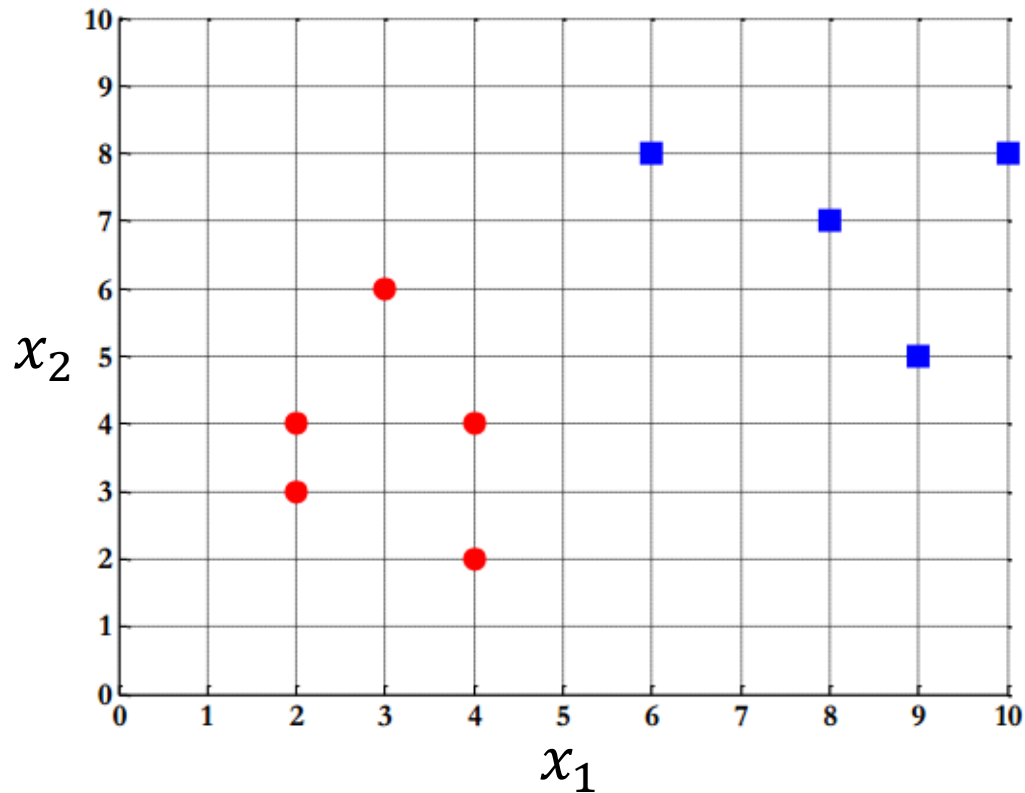
- **Workflow of LDA for the binary classification**

1. Build \mathbf{X}_1 and \mathbf{X}_2 from the training set
2. Compute $\boldsymbol{\mu}_1$ and $\boldsymbol{\mu}_2$
3. Compute \mathbf{S}_w
4. Compute \mathbf{S}_w^{-1}
5. Compute $\boldsymbol{\theta}^* = \mathbf{S}_w^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$
6. Given a testing sample, $y = \boldsymbol{\theta}^{*T} \mathbf{x}$
7. Set the threshold $\gamma = \frac{n_1 \boldsymbol{\theta}^{*T} \boldsymbol{\mu}_1 + n_2 \boldsymbol{\theta}^{*T} \boldsymbol{\mu}_2}{n_1 + n_2}$.
8. Compare y with γ to determine the class.

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

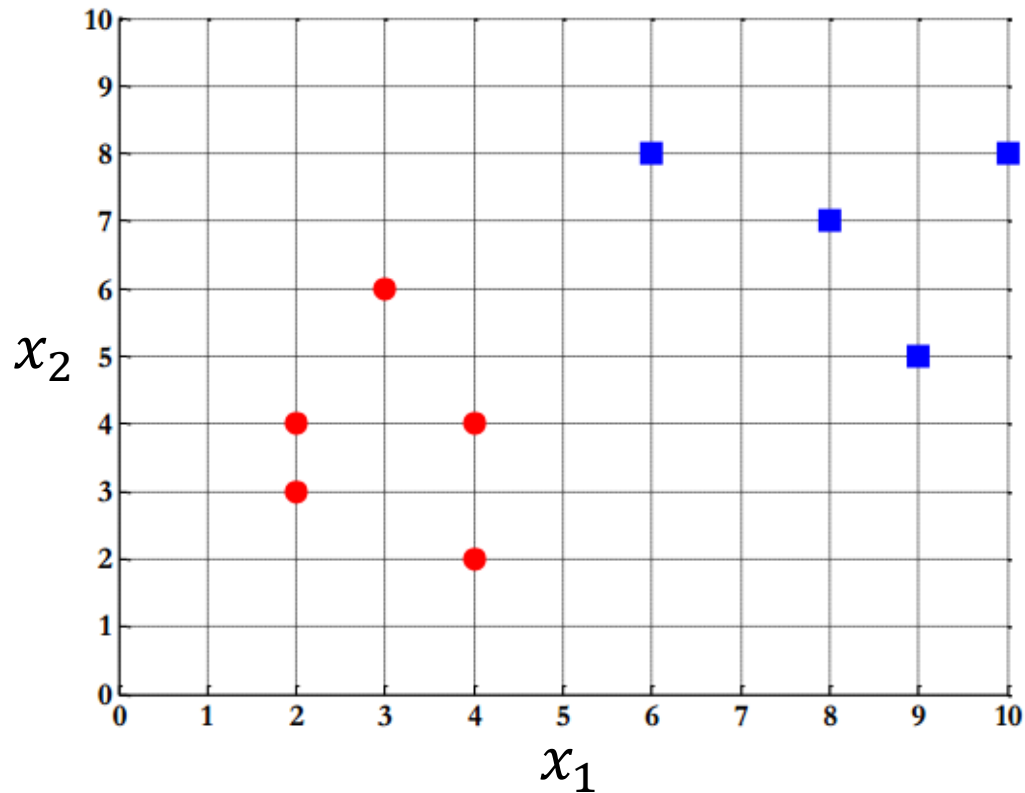
- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



- Mean of each class:

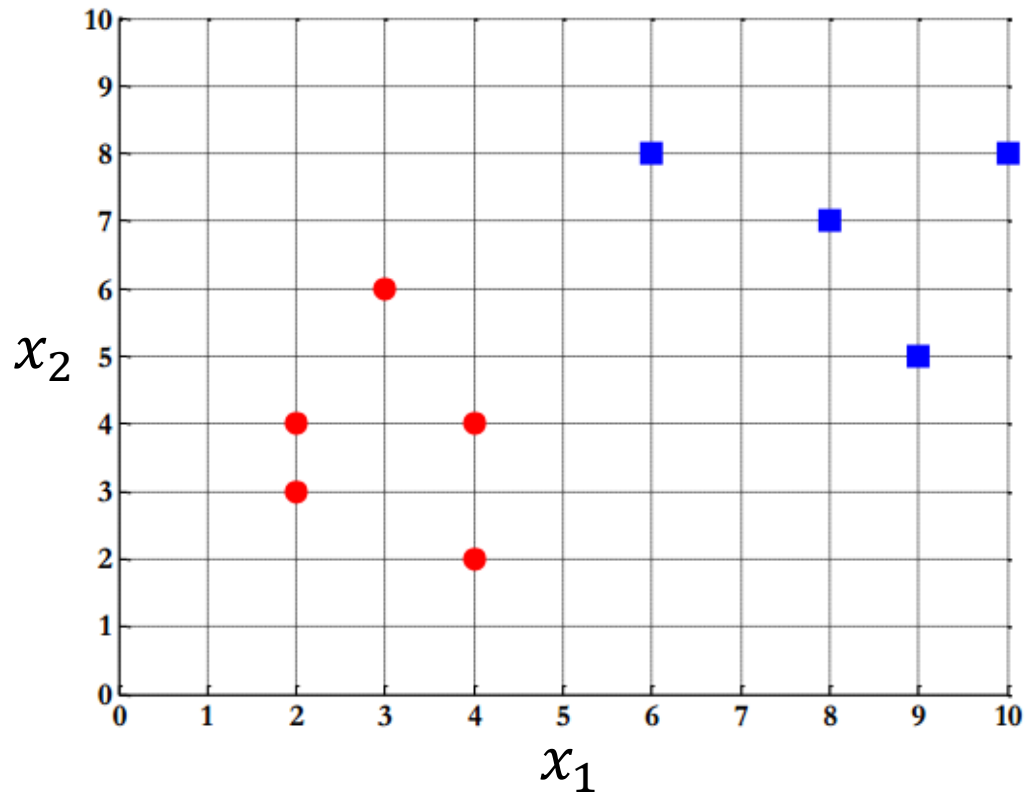
$$\mu_1 = \frac{1}{N_1} \sum_{x \in \omega_1} x = \frac{1}{5} \left[\begin{pmatrix} 4 \\ 2 \end{pmatrix} + \begin{pmatrix} 2 \\ 4 \end{pmatrix} + \begin{pmatrix} 2 \\ 3 \end{pmatrix} + \begin{pmatrix} 3 \\ 6 \end{pmatrix} + \begin{pmatrix} 4 \\ 4 \end{pmatrix} \right] = \begin{pmatrix} 3 \\ 3.8 \end{pmatrix}$$

$$\mu_2 = \frac{1}{N_2} \sum_{x \in \omega_2} x = \frac{1}{5} \left[\begin{pmatrix} 9 \\ 10 \end{pmatrix} + \begin{pmatrix} 6 \\ 8 \end{pmatrix} + \begin{pmatrix} 9 \\ 5 \end{pmatrix} + \begin{pmatrix} 8 \\ 7 \end{pmatrix} + \begin{pmatrix} 10 \\ 8 \end{pmatrix} \right] = \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix}$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



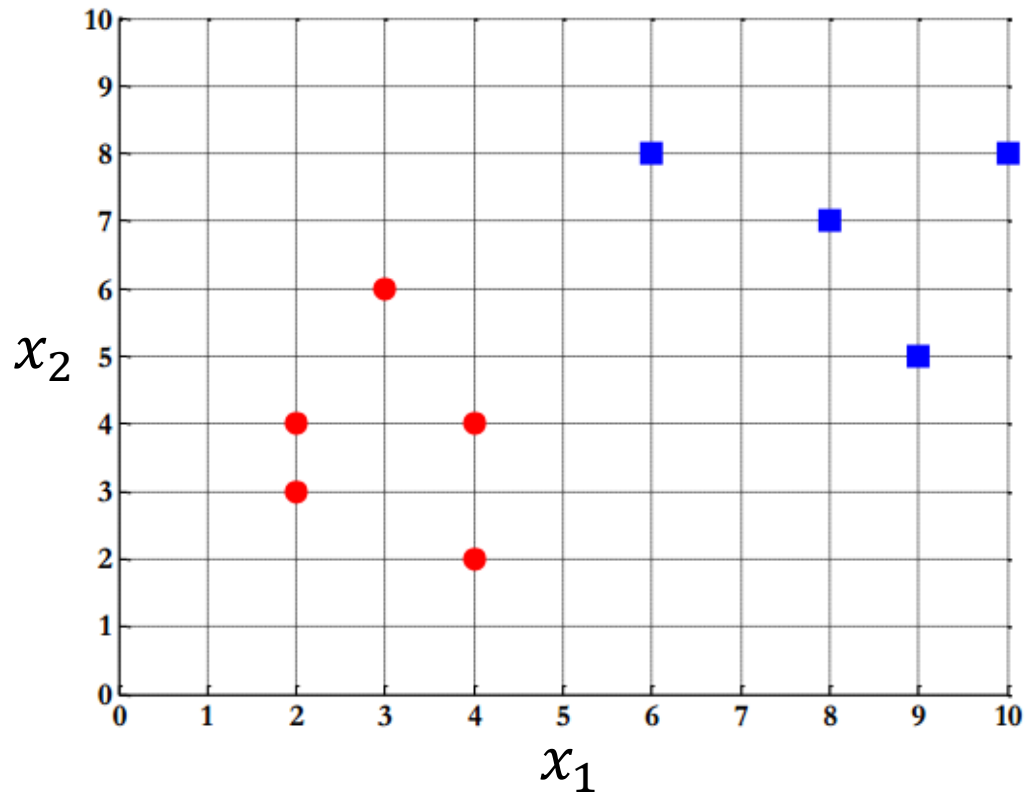
- Covariance matrix of the first class:

$$\begin{aligned} S_1 &= \sum_{x \in \omega_1} (x - \mu_1)(x - \mu_1)^T = \left[\begin{pmatrix} 4 \\ 2 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \\ &\quad + \left[\begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 3 \\ 6 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 4 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 \\ &= \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} \end{aligned}$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



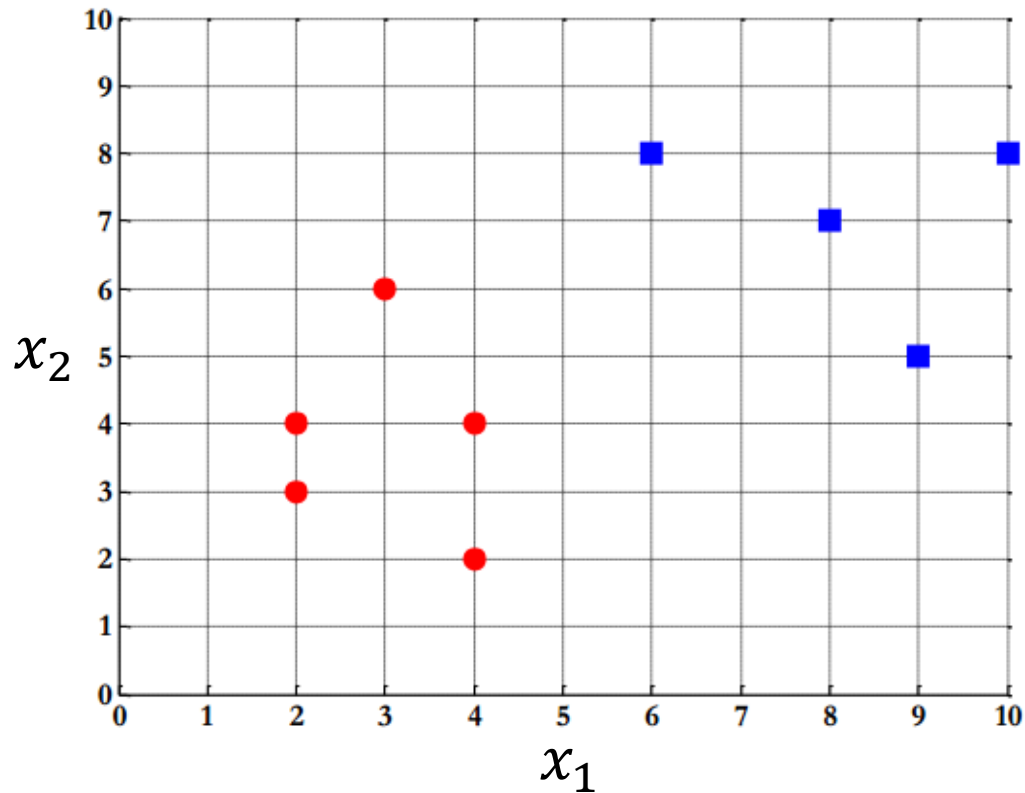
- Covariance matrix of the second class:

$$\begin{aligned} S_2 &= \sum_{x \in \omega_2} (x - \mu_2)(x - \mu_2)^T = \left[\begin{pmatrix} 9 \\ 10 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 6 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \\ &\quad + \left[\begin{pmatrix} 9 \\ 5 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 8 \\ 7 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 10 \\ 8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^2 \\ &= \begin{pmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{pmatrix} \end{aligned}$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



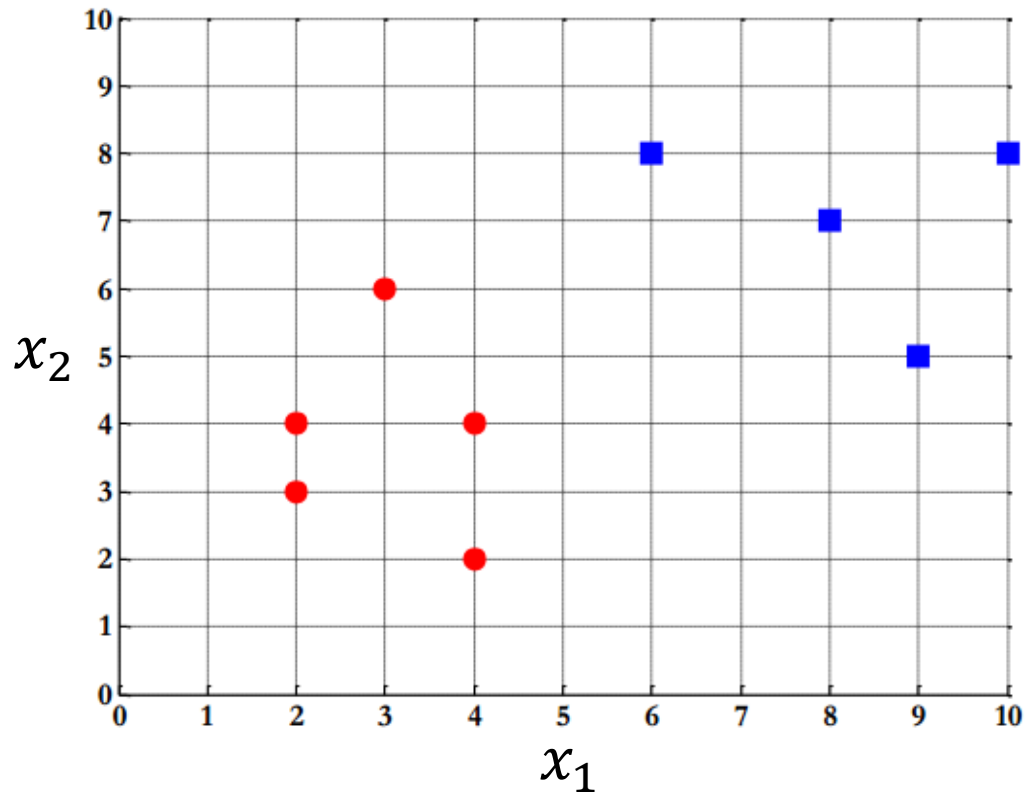
- Within-class scatter matrix:

$$\begin{aligned} S_w = S_1 + S_2 &= \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} + \begin{pmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{pmatrix} \\ &= \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix} \end{aligned}$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



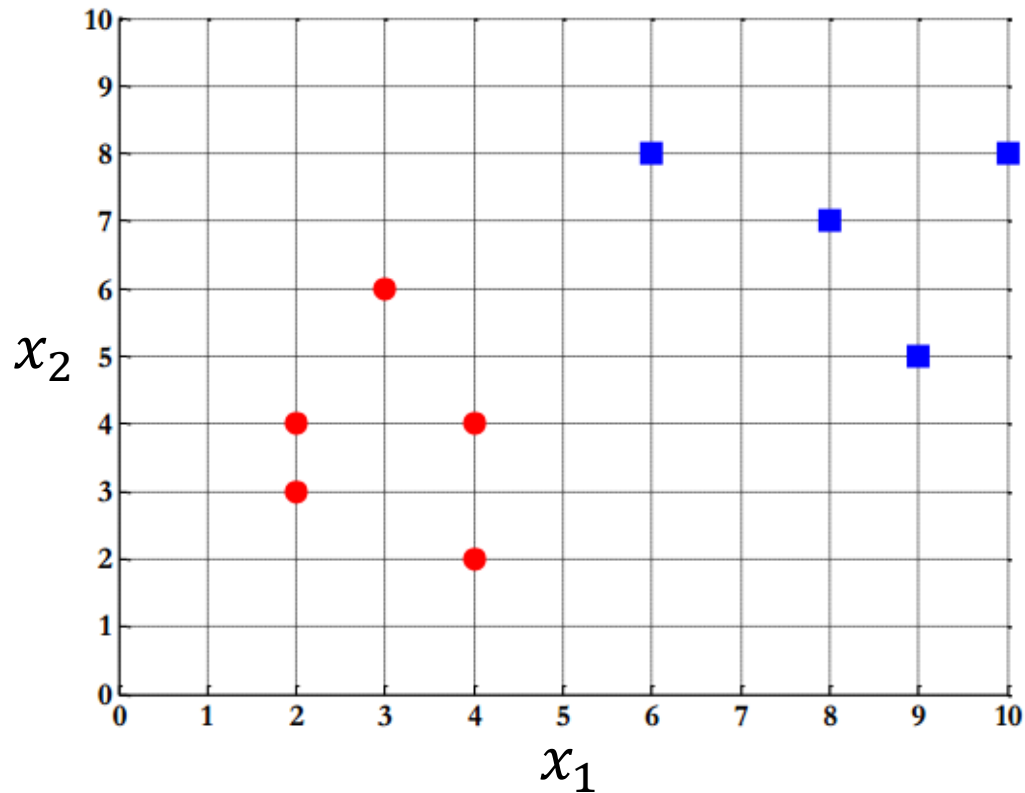
- Between-class scatter matrix:

$$\begin{aligned} S_B &= (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \\ &= \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right] \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^T \\ &= \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix} \begin{pmatrix} -5.4 & -3.8 \end{pmatrix} \\ &= \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} \end{aligned}$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



$$S_W^{-1} S_B w = \lambda w$$

$$\Rightarrow |S_W^{-1} S_B - \lambda I| = 0$$

$$\Rightarrow \begin{vmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{vmatrix}^{-1} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0$$

$$\Rightarrow \begin{vmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{vmatrix} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0$$

$$\Rightarrow \begin{vmatrix} 9.2213 - \lambda & 6.489 \\ 4.2339 & 2.9794 - \lambda \end{vmatrix}$$

$$= (9.2213 - \lambda)(2.9794 - \lambda) - 6.489 \times 4.2339 = 0$$

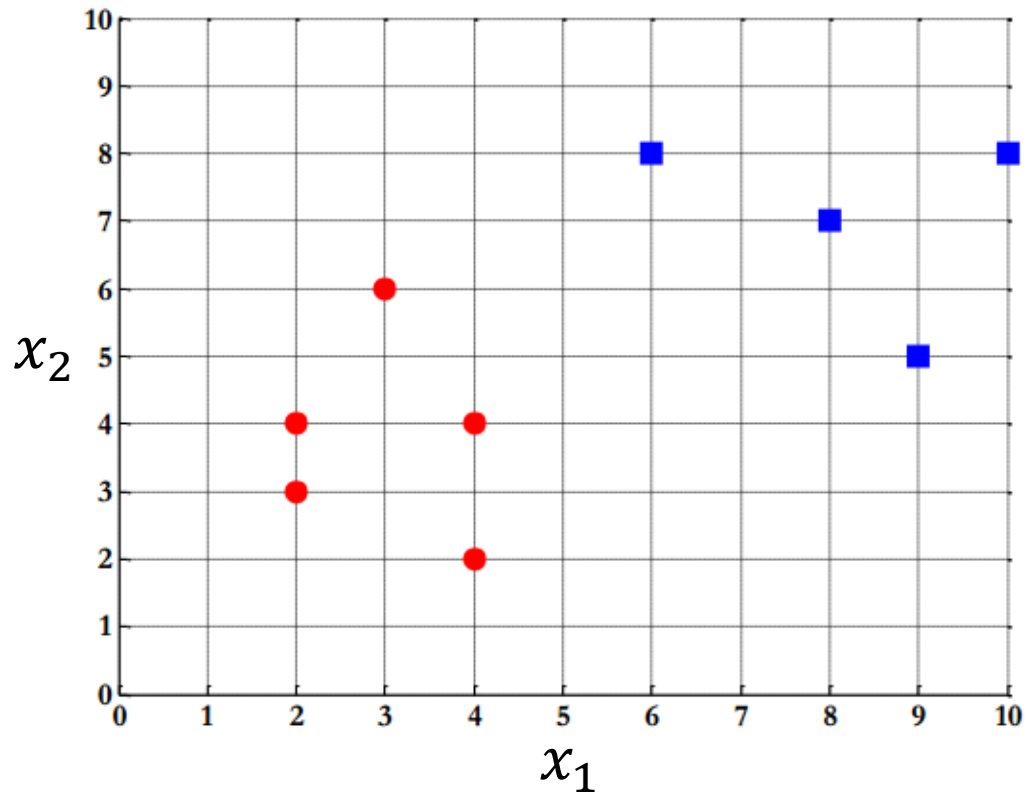
$$\Rightarrow \lambda^2 - 12.2007\lambda = 0 \Rightarrow \lambda(\lambda - 12.2007) = 0$$

$$\Rightarrow \lambda_1 = 0, \lambda_2 = 12.2007$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



$$S_W^{-1} S_B w = \lambda w$$

$$\Rightarrow |S_W^{-1} S_B - \lambda I| = 0$$

$$\Rightarrow \begin{vmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{vmatrix}^{-1} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0$$

$$\Rightarrow \begin{vmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{vmatrix} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0$$

$$\Rightarrow \begin{vmatrix} 9.2213 - \lambda & 6.489 \\ 4.2339 & 2.9794 - \lambda \end{vmatrix}$$

$$= (9.2213 - \lambda)(2.9794 - \lambda) - 6.489 \times 4.2339 = 0$$

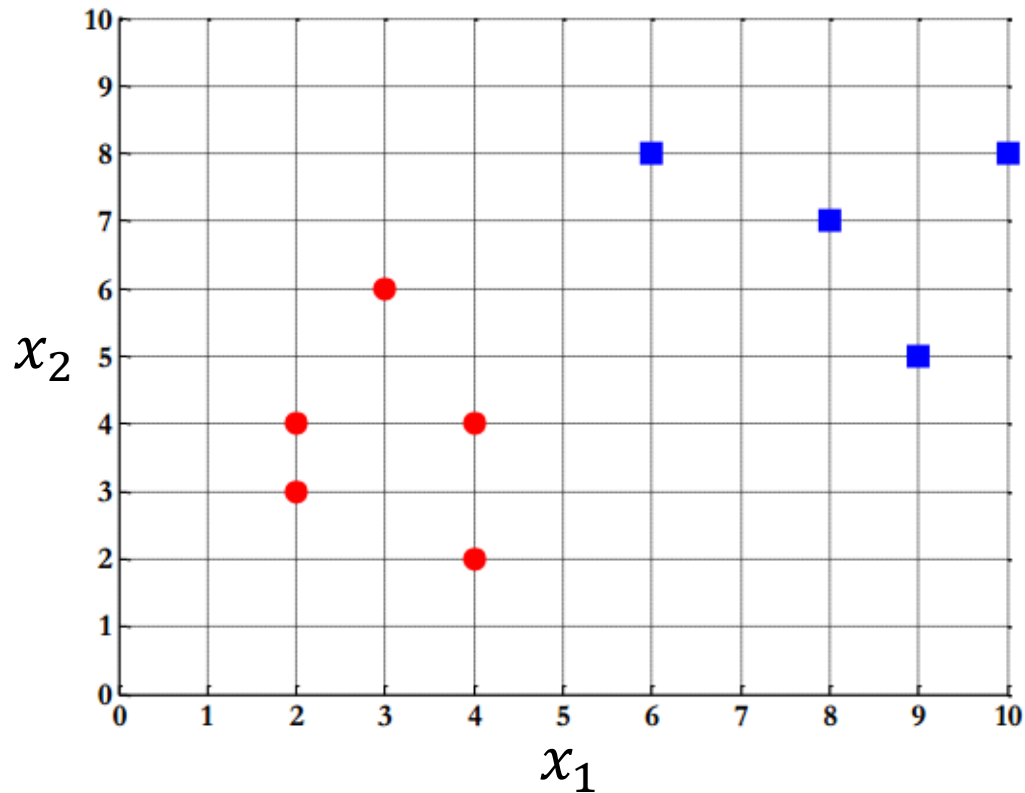
$$\Rightarrow \lambda^2 - 12.2007\lambda = 0 \Rightarrow \lambda(\lambda - 12.2007) = 0$$

$$\Rightarrow \lambda_1 = 0, \lambda_2 = 12.2007$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class $\omega_1: \mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class $\omega_2: \mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



- The optimal projection is the one that given maximum $J = -\boldsymbol{\theta}^T \mathbf{S}_b \boldsymbol{\theta} = -\lambda$

$$\begin{pmatrix} 9.2213 & 6.489 \\ 4.2339 & 2.9794 \end{pmatrix} w_1 = 0 \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

and

$$\begin{pmatrix} 9.2213 & 6.489 \\ 4.2339 & 2.9794 \end{pmatrix} w_2 = \underbrace{12.2007}_{\lambda_2} \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$$

Thus;

$$w_1 = \begin{pmatrix} -0.5755 \\ 0.8178 \end{pmatrix}$$

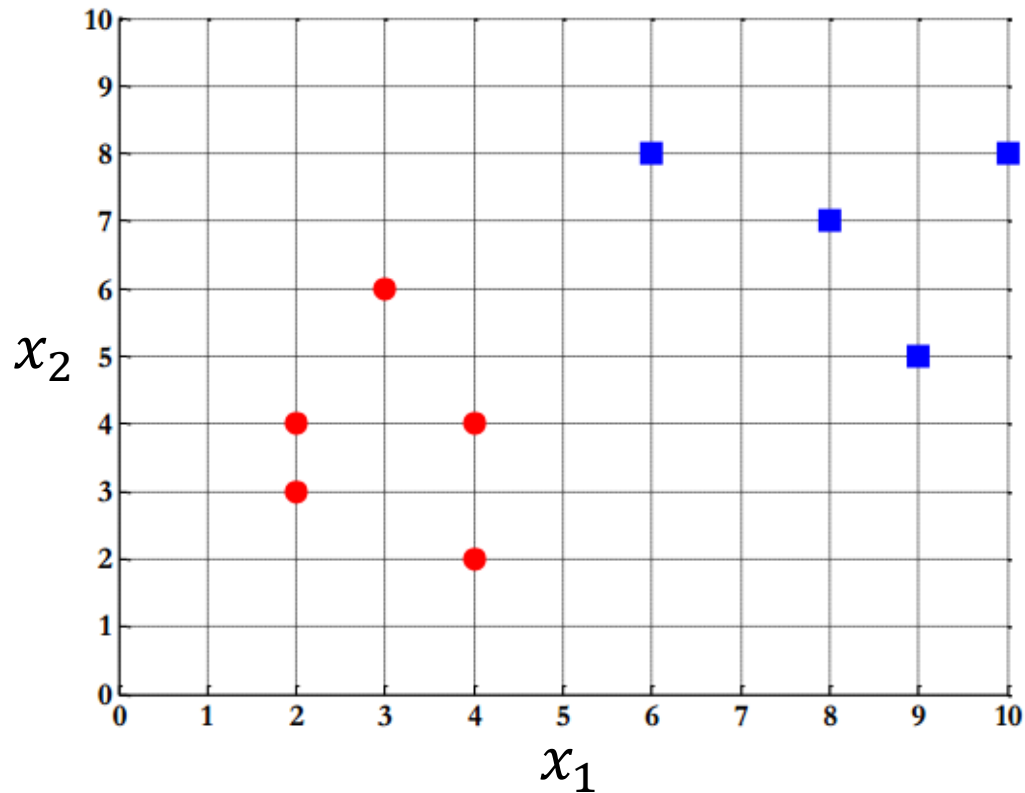
and

$$w_2 = \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix} = w^*$$

Example

Compute the Linear Discriminant projection for the following two dimensional dataset.

- Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$
- Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$



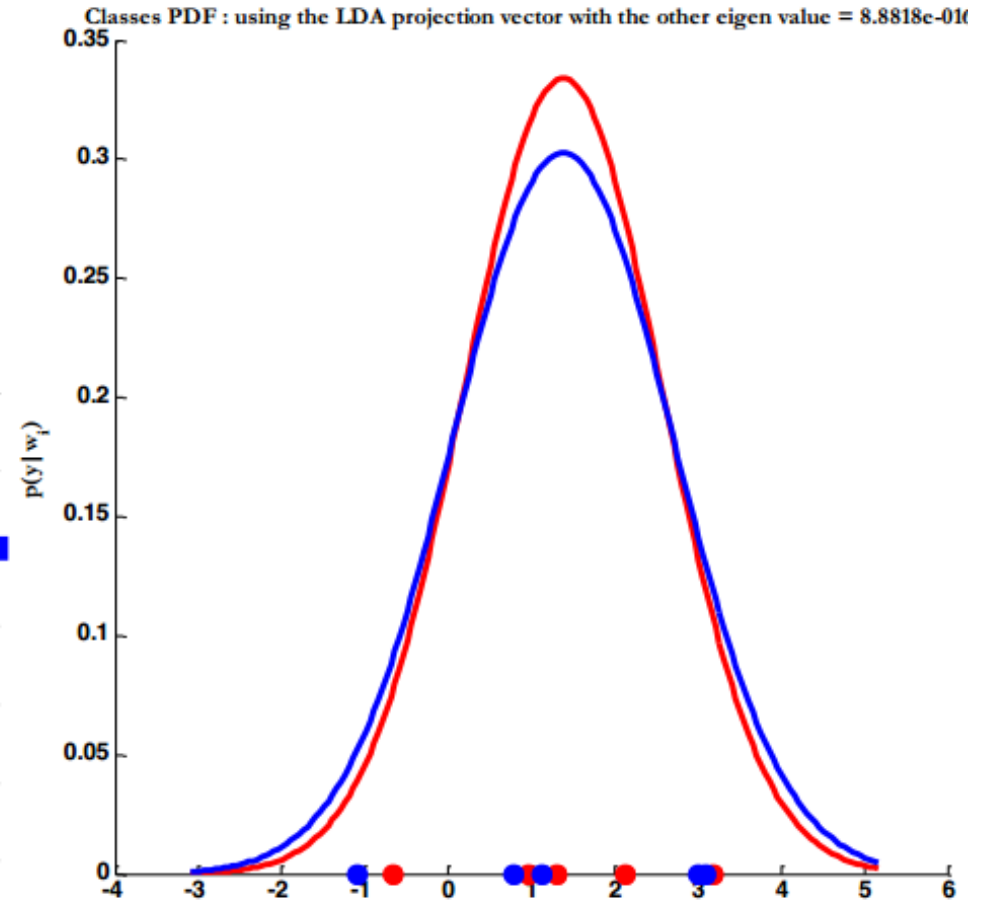
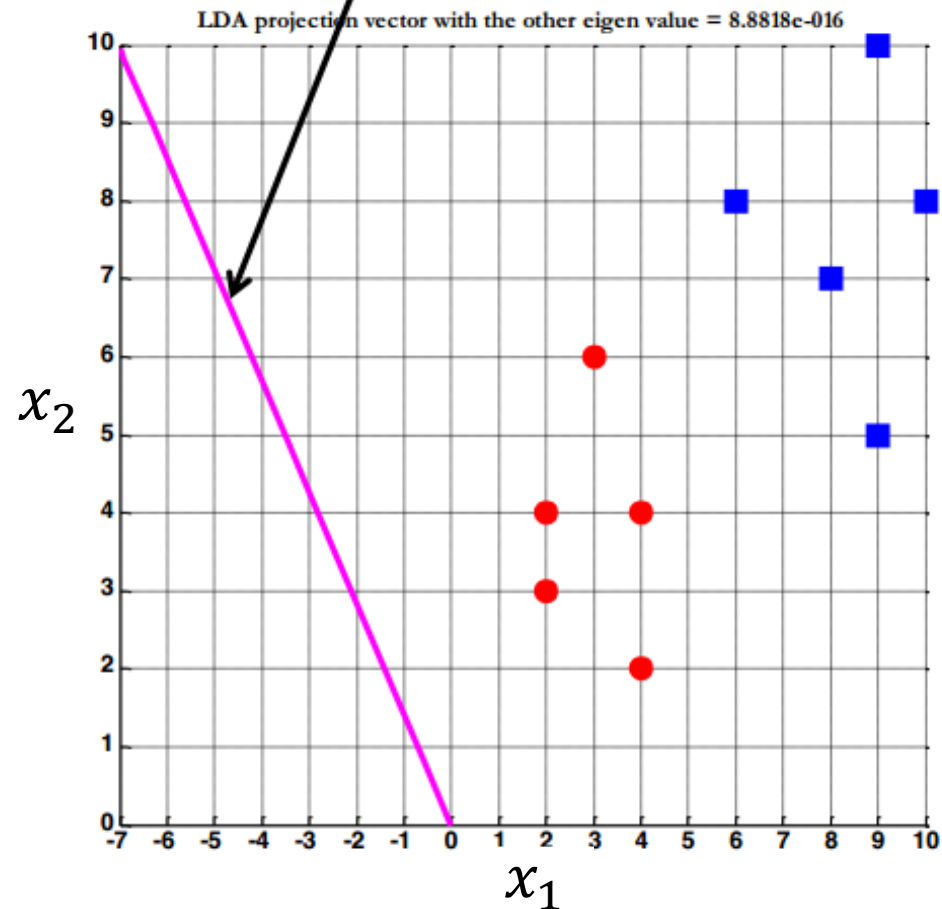
- Or directly,

$$\begin{aligned} w^* &= S_W^{-1}(\mu_1 - \mu_2) = \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix}^{-1} \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right] \\ &= \begin{pmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{pmatrix} \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix} \\ &= \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix} \end{aligned}$$

Example

LDA--Projection

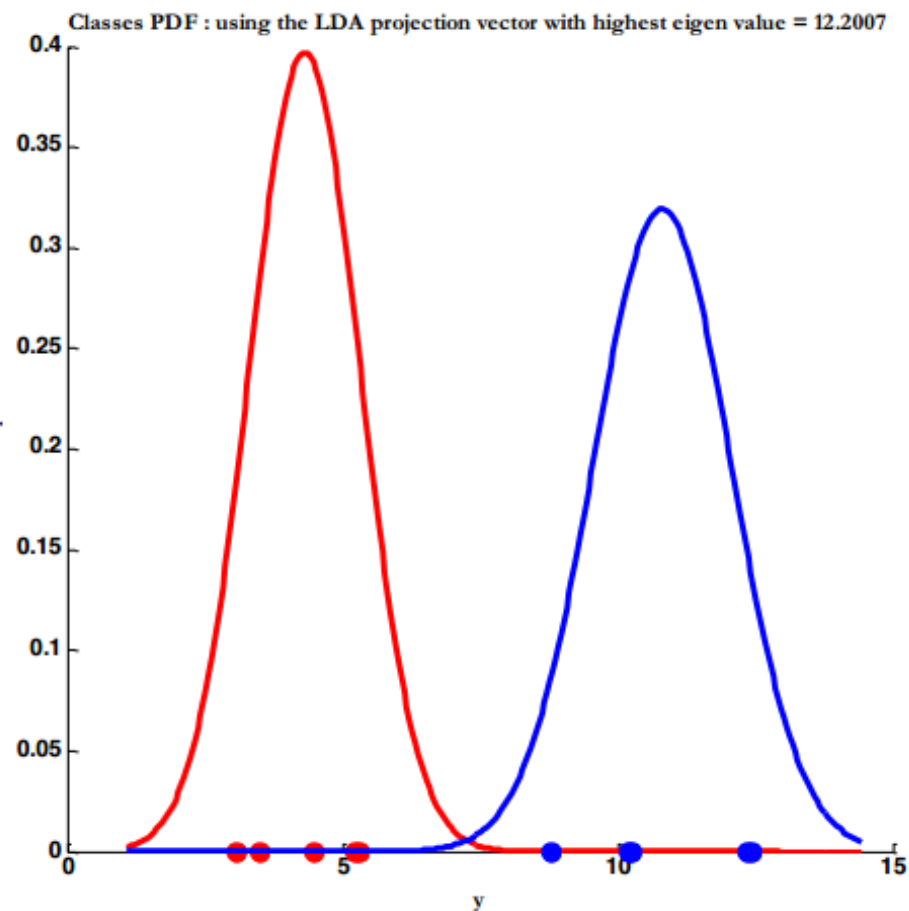
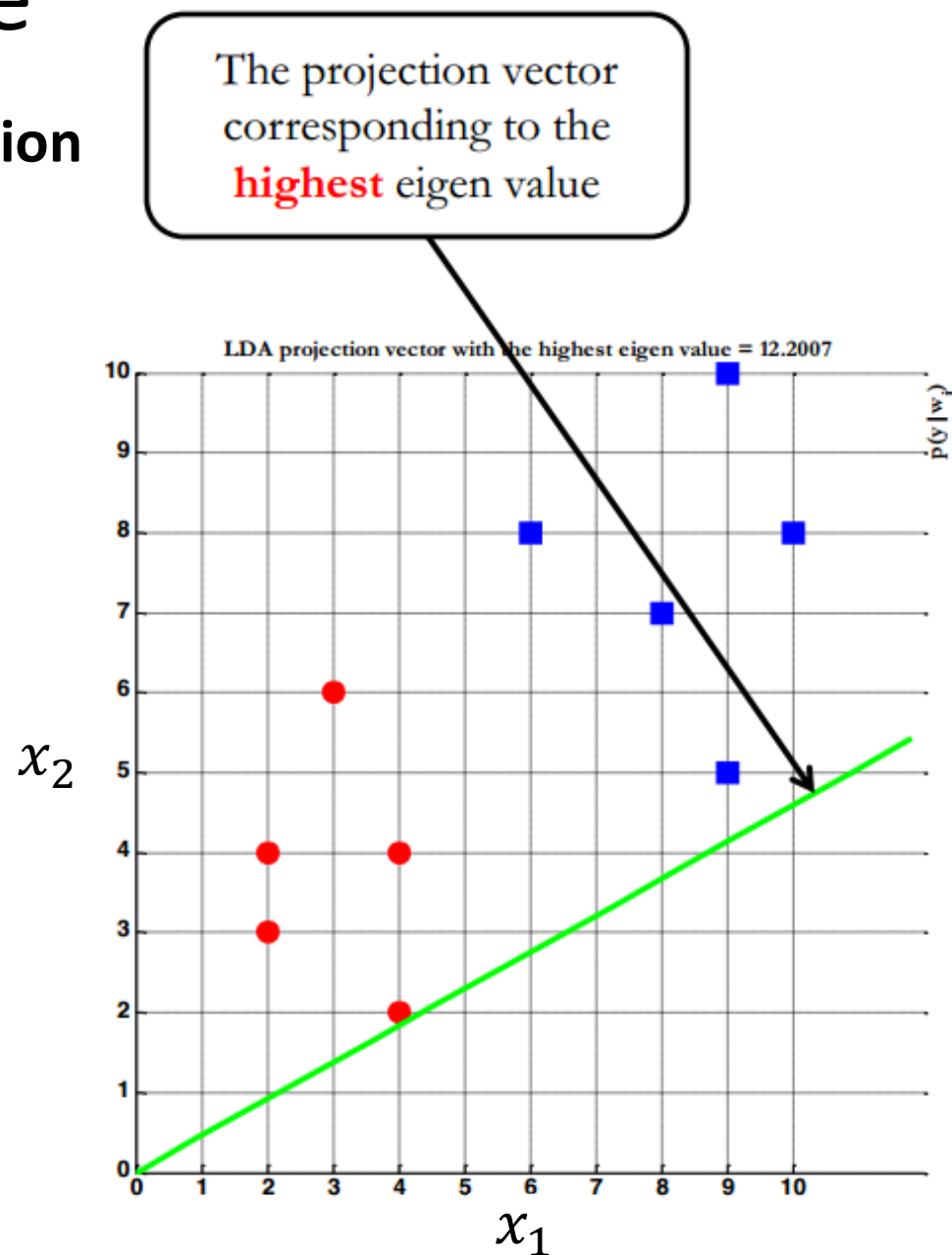
The projection vector
corresponding to the
smallest eigen value



Using this vector leads to
bad separability
between the two classes

Example

LDA--Projection



Using this vector leads to **good separability** between the two classes

Linear Discriminant Analysis— C Classes

- Assume we have C classes, each class has n_i d -dimensional samples, where $i = 1, 2, \dots, C$
- A transformation $\Theta \in \mathbb{R}^{d \times p}$: project the samples in \mathbf{X} onto \mathbf{Y} ($p \ll d$). In fact, $p \leq C + 1$, we will see later.

$$\mathbf{y}_i = \Theta^T \mathbf{x}_i$$

$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix}$$

$$\mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{ip} \end{bmatrix}$$

$$\Theta = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_p] \in \mathbb{R}^{d \times p}$$

Statistical Facts

Class mean vector (sample):

$$\boldsymbol{\mu}_i = \frac{1}{n_i} \sum_{\mathbf{x} \in C_i} \mathbf{x}, \boldsymbol{\mu}_i \in \mathbb{R}^{d \times 1}$$

Within-class scatter:

$$\mathbf{S}_w = \sum_{i=1}^C \mathbf{S}_{wi} \quad \mathbf{S}_{wi} = \sum_{\mathbf{x} \in C_i} (\mathbf{x} - \boldsymbol{\mu}_i)(\mathbf{x} - \boldsymbol{\mu}_i)^T \quad \mathbf{S}_w \in \mathbb{R}^{d \times d}$$

Between-class scatter:

$$\mathbf{S}_b = \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T = \frac{1}{2N} \sum_{i,j=1}^C n_i n_j (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T \quad \mathbf{S}_b \in \mathbb{R}^{d \times d}$$

Statistical Facts

Total covariance (sample):

$$\mathbf{S}_t = \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \mathbf{S}_w + \mathbf{S}_b$$

$$\begin{aligned}\mathbf{S}_t &= \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \sum_{i=1}^C \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu})^T & \mathbf{x}_{ij} \in C_i \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i) + (\boldsymbol{\mu}_i - \boldsymbol{\mu})][(\mathbf{x}_{ij} - \boldsymbol{\mu}_i) + (\boldsymbol{\mu}_i - \boldsymbol{\mu})]^T \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T] \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T] = \mathbf{S}_w + \mathbf{S}_b\end{aligned}$$

Statistical Facts

Total covariance (sample):

$$\mathbf{S}_t = \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \mathbf{S}_w + \mathbf{S}_b$$

$$\begin{aligned}\mathbf{S}_t &= \sum_{\mathbf{x}} (\mathbf{x} - \boldsymbol{\mu})(\mathbf{x} - \boldsymbol{\mu})^T = \sum_{i=1}^C \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu})^T & \mathbf{x}_{ij} \in C_i \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i) + (\boldsymbol{\mu}_i - \boldsymbol{\mu})][(\mathbf{x}_{ij} - \boldsymbol{\mu}_i) + (\boldsymbol{\mu}_i - \boldsymbol{\mu})]^T \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T] \\ &= \sum_{i=1}^C \sum_{j=1}^{n_i} [(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T] = \mathbf{S}_w + \mathbf{S}_b\end{aligned}$$

$$\sum_{i=1}^C \sum_{j=1}^{n_i} (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\mathbf{x}_{ij} - \boldsymbol{\mu}_i)^T = \sum_{i=1}^C (\boldsymbol{\mu}_i - \boldsymbol{\mu}) \left(\sum_{j=1}^{n_i} \mathbf{x}_{ij} - \sum_{j=1}^{n_i} \boldsymbol{\mu}_i \right)^T = 0$$

Linear Discriminant Analysis— C Classes

- Assume we have C classes, each class has n_i d -dimensional samples, where $i = 1, 2, \dots, C$
- A transformation $\Theta \in \mathbb{R}^{d \times p}$: project the samples in \mathbf{X} onto \mathbf{Y} ($p \ll d$). In fact, $p \leq C + 1$, we will see later.

$$\mathbf{y}_i = \Theta^T \mathbf{x}_i$$
$$\mathbf{x}_i = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{id} \end{bmatrix} \quad \mathbf{y}_i = \begin{bmatrix} y_{i1} \\ y_{i2} \\ \vdots \\ y_{ip} \end{bmatrix} \quad \Theta = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_p] \in \mathbb{R}^{d \times p}$$

$$\tilde{\mathbf{S}}_w = \Theta^T \mathbf{S}_w \Theta \quad \tilde{\mathbf{S}}_b = \Theta^T \mathbf{S}_b \Theta \quad \tilde{\boldsymbol{\mu}}_i = \Theta^T \boldsymbol{\mu}_i \quad \tilde{\boldsymbol{\mu}} = \Theta^T \boldsymbol{\mu}$$

Linear Discriminant Analysis— C Classes

Popular objective function:

$$J_1(\Theta) = \max_{\Theta} \frac{\text{tr}(\tilde{\mathbf{S}}_b)}{\text{tr}(\tilde{\mathbf{S}}_w)} = \max_{\Theta} \frac{\text{tr}(\Theta^T \mathbf{S}_b \Theta)}{\text{tr}(\Theta^T \mathbf{S}_w \Theta)}$$

$$J_2(\Theta) = \max_{\Theta} \text{tr}(\tilde{\mathbf{S}}_w^{-1} \tilde{\mathbf{S}}_b) = \max_{\Theta} \text{tr}((\Theta^T \mathbf{S}_w \Theta)^{-1} \Theta^T \mathbf{S}_b \Theta)$$

$$J_3(\Theta) = \frac{|\tilde{\mathbf{S}}_b|}{|\tilde{\mathbf{S}}_w|}$$

This technique was developed by R. A. Fisher (1936) for **the two-class case** and extended by C. R. Rao (1948) to handle **the multiclass case**.

Linear Discriminant Analysis— C Classes

In $J_1(\mathbf{\Theta})$, what is the meaning of “**trace**”?

$$J_1(\mathbf{\Theta}) = \max_{\mathbf{\Theta}} \frac{\text{tr}(\tilde{\mathbf{S}}_b)}{\text{tr}(\tilde{\mathbf{S}}_w)} = \max_{\mathbf{\Theta}} \frac{\text{tr}(\mathbf{\Theta}^T \mathbf{S}_b \mathbf{\Theta})}{\text{tr}(\mathbf{\Theta}^T \mathbf{S}_w \mathbf{\Theta})}$$

$$\mathbf{\Theta}^T \mathbf{S}_b \mathbf{\Theta} = \begin{bmatrix} \boldsymbol{\theta}_1^T \\ \vdots \\ \boldsymbol{\theta}_p^T \end{bmatrix} \mathbf{S}_b [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_p] = \begin{bmatrix} \boldsymbol{\theta}_1^T \\ \vdots \\ \boldsymbol{\theta}_p^T \end{bmatrix} [\mathbf{S}_b \boldsymbol{\theta}_1, \mathbf{S}_b \boldsymbol{\theta}_2, \dots, \mathbf{S}_b \boldsymbol{\theta}_p]$$

Linear Discriminant Analysis— C Classes

In $J_1(\Theta)$, what is the meaning of “**trace**”?

$$J_1(\Theta) = \max_{\Theta} \frac{\text{tr}(\tilde{\mathbf{S}}_b)}{\text{tr}(\tilde{\mathbf{S}}_w)} = \max_{\Theta} \frac{\text{tr}(\Theta^T \mathbf{S}_b \Theta)}{\text{tr}(\Theta^T \mathbf{S}_w \Theta)}$$

$$\Theta^T \mathbf{S}_b \Theta = \begin{bmatrix} \theta_1^T \\ \vdots \\ \theta_p^T \end{bmatrix} \mathbf{S}_b [\theta_1, \theta_2, \dots, \theta_p] = \begin{bmatrix} \theta_1^T \\ \vdots \\ \theta_p^T \end{bmatrix} [\mathbf{S}_b \theta_1, \mathbf{S}_b \theta_2, \dots, \mathbf{S}_b \theta_p]$$

$$\text{tr}(\Theta^T \mathbf{S}_b \Theta) = \sum_{i=1}^p \theta_i^T \mathbf{S}_b \theta_i \qquad \text{tr}(\Theta^T \mathbf{S}_w \Theta) = \sum_{i=1}^p \theta_i^T \mathbf{S}_w \theta_i$$

Linear Discriminant Analysis— C Classes

Optimization $J_1(\Theta)$:

- Recall in two-classes case, we solved the eigen value problem.

$$\begin{array}{ll} \min_{\theta} -\theta^T S_b \theta & \Rightarrow S_b \theta = \lambda S_w \theta \\ \text{s.t. } \theta^T S_w \theta = 1 & \end{array}$$

Linear Discriminant Analysis— C Classes

Optimization $J_1(\Theta)$:

- Recall in two-classes case, we solved the eigen value problem.

$$\begin{array}{ll} \min_{\theta} -\theta^T \mathbf{S}_b \theta & \Rightarrow \quad \mathbf{S}_b \theta = \lambda \mathbf{S}_w \theta \\ \text{s.t. } \theta^T \mathbf{S}_w \theta = 1 & \end{array}$$

- For C -classes case, we have p projection vectors,

$$\mathbf{S}_w^{-1} \mathbf{S}_b \theta_i = \lambda \theta_i, \quad i = 1, 2, \dots, p$$

Linear Discriminant Analysis— C Classes

Optimization $J_1(\Theta)$:

- Recall in two-classes case, we solved the eigen value problem.

$$\begin{aligned} \min_{\theta} -\theta^T S_b \theta \\ \text{s.t. } \theta^T S_w \theta = 1 \end{aligned} \quad \Rightarrow \quad S_b \theta = \lambda S_w \theta$$

- For C -classes case, we have p projection vectors,

$$S_w^{-1} S_b \theta_i = \lambda \theta_i, \quad i = 1, 2, \dots, p$$

Columns of Θ^* are eigenvectors corresponding to the largest eigenvalues:

$$S_w^{-1} S_b \Theta^* = \lambda \Theta^* \quad \Theta^* = [\theta_1^*, \theta_2^*, \dots, \theta_p^*]$$

Linear Discriminant Analysis— C Classes

Optimization $J_1(\Theta)$:

- Recall in two-classes case, we solved the eigen value problem.

$$\begin{aligned} \min_{\theta} -\theta^T S_b \theta \\ \text{s.t. } \theta^T S_w \theta = 1 \end{aligned} \quad \Rightarrow \quad S_b \theta = \lambda S_w \theta$$

- For C -classes case, we have p projection vectors,

$$S_w^{-1} S_b \theta_i = \lambda \theta_i, \quad i = 1, 2, \dots, p$$

Columns of Θ^* are eigenvectors corresponding to the largest eigenvalues:

$$S_w^{-1} S_b \Theta^* = \lambda \Theta^* \quad \Theta^* = [\theta_1^*, \theta_2^*, \dots, \theta_p^*]$$

$p \leq C + 1$, why?

Linear Discriminant Analysis— C Classes

Optimization $J_1(\Theta)$:

- \mathbf{S}_b has a maximum rank of $C - 1$.
- \mathbf{S}_b is the sum of C *rank* = 1 matrices, and because only $C - 1$ of these are independent,

$$\mathbf{S}_b = \sum_{i=1}^C \frac{n_i}{N} (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T$$

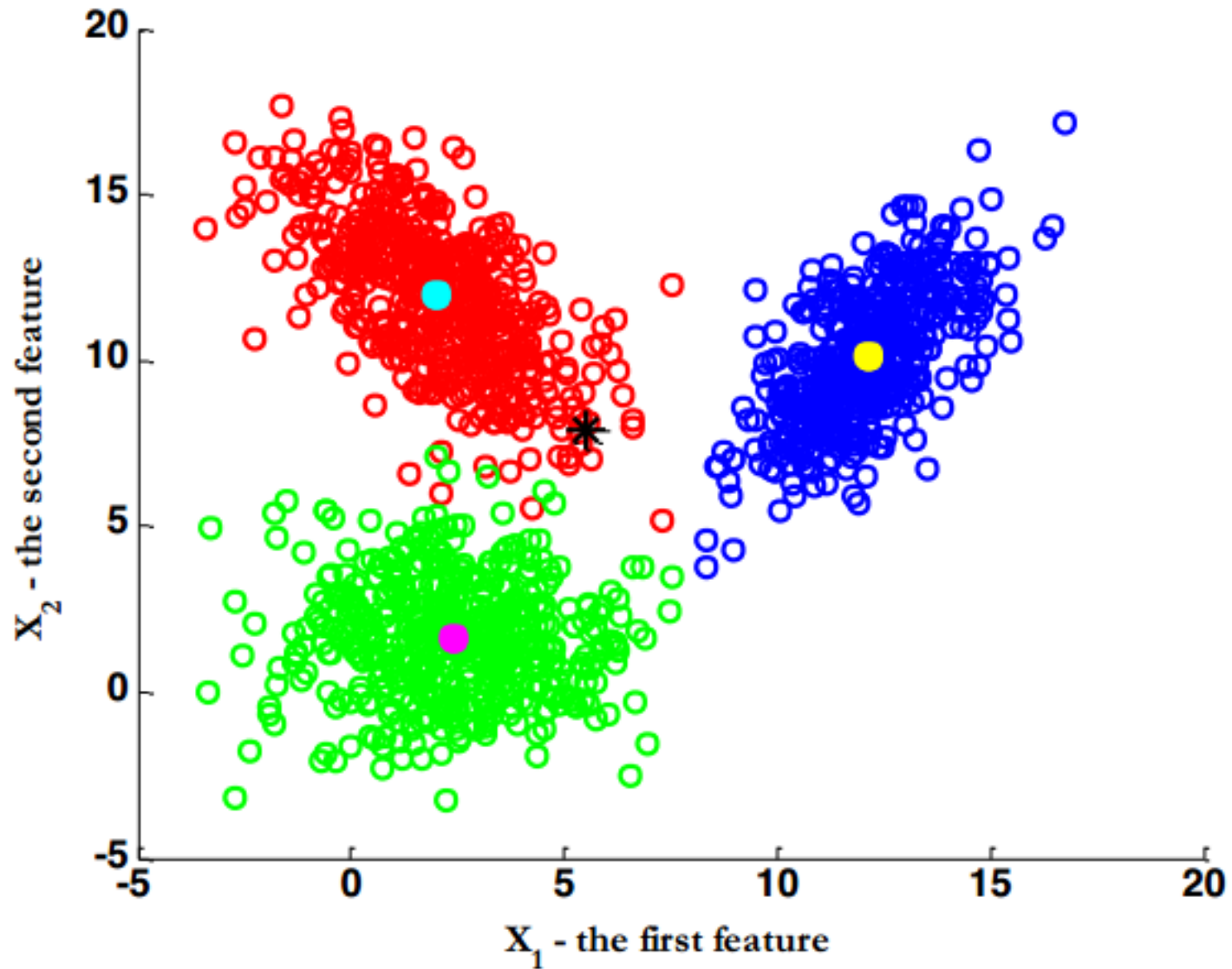
- Given a matrix $\mathbf{A}_{m \times n}$ and $\mathbf{B}_{n \times k}$,
- $\text{rank}(\mathbf{AB}) = \min(\text{rank}(\mathbf{A}), \text{rank}(\mathbf{B}))$
- $\text{rank}(\mathbf{A} + \mathbf{B}) \leq \text{rank}(\mathbf{A}) + \text{rank}(\mathbf{B})$

$$\text{rank}\left((\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T\right) \leq \text{rank}(\boldsymbol{\mu}_i - \boldsymbol{\mu}) = 1 \qquad \text{rank}(\mathbf{S}_w^{-1} \mathbf{S}_b) \leq \text{rank}(\mathbf{S}_b) \leq C - 1$$

Linear Discriminant Analysis— C Classes

- **Workflow of LDA for the C -classification**
 1. Compute $\boldsymbol{\mu}_i$
 2. Compute \boldsymbol{S}_b
 3. Compute \boldsymbol{S}_w^{-1}
 4. Compute the largest p eigenvalues of $\boldsymbol{S}_w^{-1}\boldsymbol{S}_b$ and the corresponding eigenvectors $\{\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_p\}$.
 5. Let $\boldsymbol{\Theta} = [\boldsymbol{\theta}_1, \boldsymbol{\theta}_2, \dots, \boldsymbol{\theta}_p]$, then $\boldsymbol{y}_i = \boldsymbol{\Theta}^T \boldsymbol{x}_i$

Illustration-3 Classes



```

%% computing the LDA
% class means
Mu1 = mean(X1')';
Mu2 = mean(X2')';
Mu3 = mean(X3')';

% overall mean
Mu = (Mu1 + Mu2 + Mu3) ./ 3;

% class covariance matrices
S1 = cov(X1');
S2 = cov(X2');
S3 = cov(X3');

% within-class scatter matrix
Sw = S1 + S2 + S3;

% number of samples of each class
N1 = size(X1,2);
N2 = size(X2,2);
N3 = size(X3,2);

% between-class scatter matrix
SB1 = N1 .* (Mu1-Mu)*(Mu1-Mu)';
SB2 = N2 .* (Mu2-Mu)*(Mu2-Mu)';
SB3 = N3 .* (Mu3-Mu)*(Mu3-Mu)';

SB = SB1 + SB2 + SB3;

% computing the LDA projection
invSw = inv(Sw);
invSw_by_SB = invSw * SB;

% getting the projection vectors
% [V,D] = EIG(X) produces a diagonal matrix D of eigenvalues and a
% full matrix V whose columns are the corresponding eigenvectors
[V,D] = eig(invSw_by_SB);

% the projection vectors - we will have at most C-1 projection vectors,
% from which we can choose the most important ones ranked by their
% corresponding eigen values ... lets investigate the two projection
% vectors
W1 = V(:,1);
W2 = V(:,2);

```

Recall ...

$$S_W = \sum_{i=1}^C S_i$$

where $S_i = \sum_{x \in \omega_i} (x - \mu_i)(x - \mu_i)^T$

and $\mu_i = \frac{1}{N_i} \sum_{x \in \omega_i} x$

$$S_B = \sum_{i=1}^C N_i (\mu_i - \mu)(\mu_i - \mu)^T$$

where $\mu = \frac{1}{N} \sum_{\forall x} x = \frac{1}{N} \sum_{\forall x} N_i \mu_i$

and $\mu_i = \frac{1}{N_i} \sum_{x \in \omega_i} x$

$$S_W^{-1} S_B$$

```

%% lets visualize them ...
% we will plot the scatter plot to better visualize the features
hfig = figure;
axes1 = axes('Parent',hfig,'FontWeight','bold','FontSize',12);
hold('all');

% Create xlabel
xlabel('X_1 - the first feature','FontWeight','bold','FontSize',12,...
    'FontName','Garamond');

% Create ylabel
ylabel('X_2 - the second feature','FontWeight','bold','FontSize',12,...
    'FontName','Garamond');

% the first class
scatter(X1(1,:),X1(2,:), 'r','LineWidth',2,'Parent',axes1);
hold on

% class's mean
plot(Mu1_est(1),Mu1_est(2),'co','MarkerSize',8,'MarkerEdgeColor','c',...
    'Color','c','LineWidth',2,'MarkerFaceColor','c','Parent',axes1);
hold on

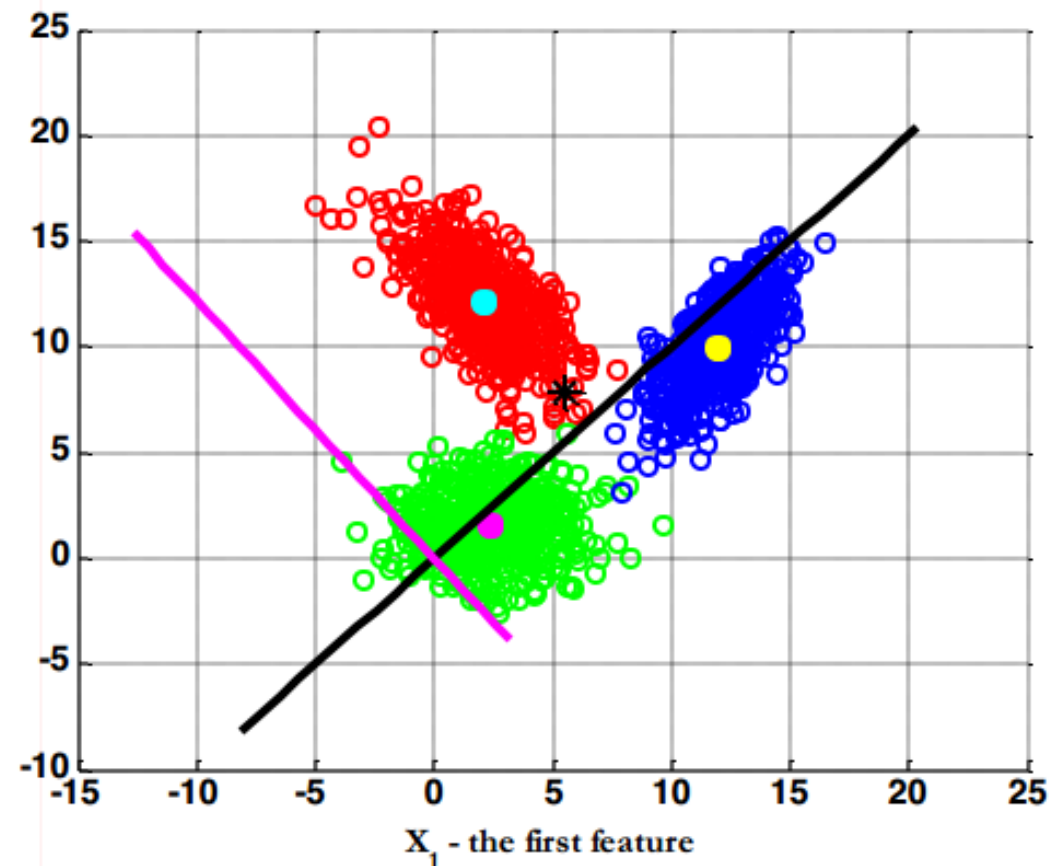
% the second class
scatter(X2(1,:),X2(2,:), 'g','LineWidth',2,'Parent',axes1);
hold on

% class's mean
plot(Mu2_est(1),Mu2_est(2),'mo','MarkerSize',8,'MarkerEdgeColor','m',...
    'Color','m','LineWidth',2,'MarkerFaceColor','m','Parent',axes1);
hold on

% the third class
scatter(X3(1,:),X3(2,:), 'b','LineWidth',2,'Parent',axes1);
hold on

% class's mean
plot(Mu3_est(1),Mu3_est(2),'yo','LineWidth',2,'MarkerSize',8,'MarkerEdgeColor',...
    'y','Color','y','MarkerFaceColor','y','Parent',axes1);
hold on

```



```

% drawing the projection vectors
% the first vector
t = -10:25;
line_x1 = t .* W1(1);
line_y1 = t .* W1(2);

% the second vector
t = -5:20;
line_x2 = t .* W2(1);
line_y2 = t .* W2(2);

plot(line_x1,line_y1,'k-', 'LineWidth', 3);
hold on
plot(line_x2,line_y2,'m-', 'LineWidth', 3);
grid on

```

Along first projection vector $y = w_1^T x$

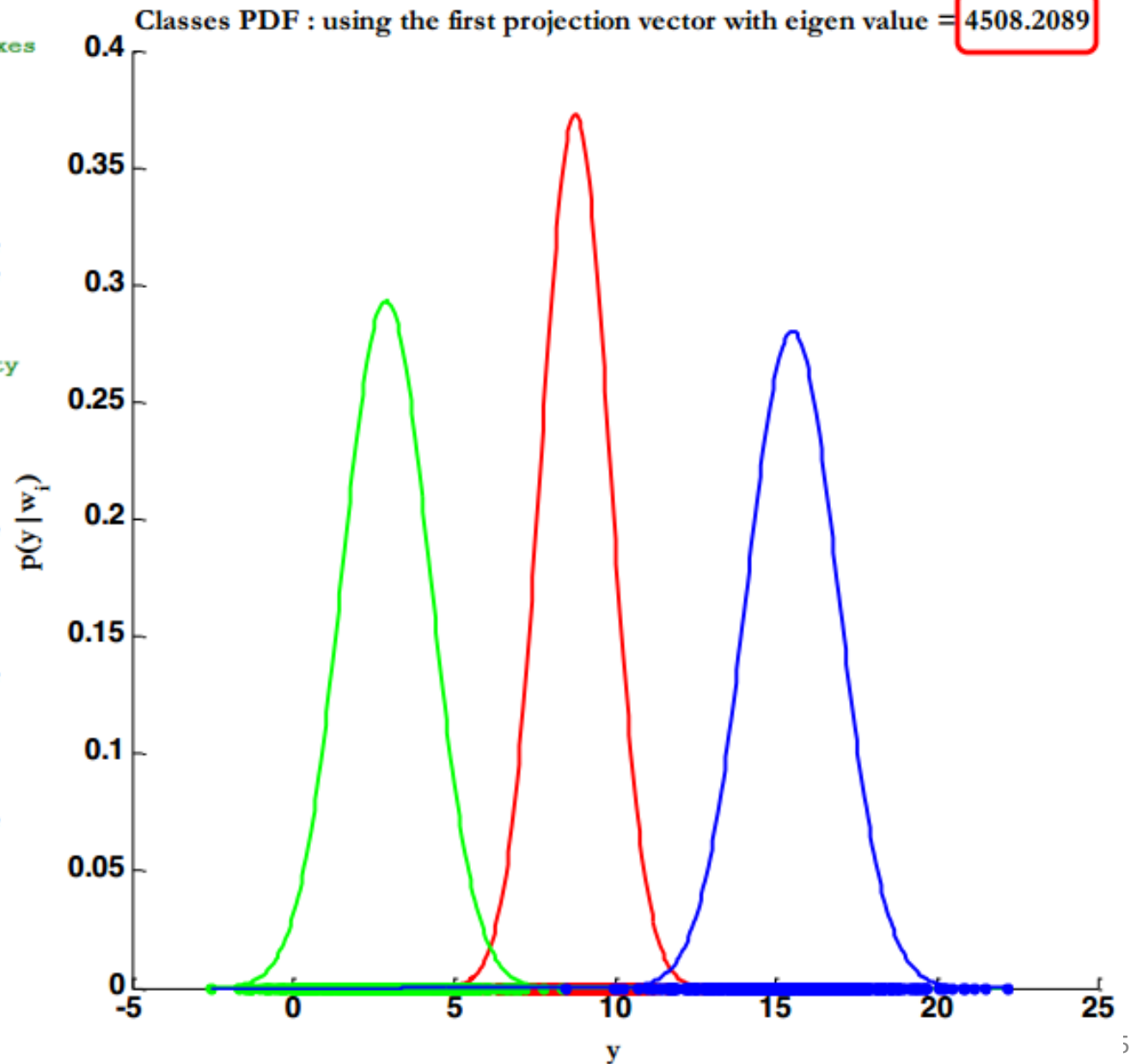
```
% project data samples along the projections axes
% the first projection vector
y1_w1 = W1'*X1;
y2_w1 = W1'*X2;
y3_w1 = W1'*X3;

% projection limits
minY = min([min(y1_w1),min(y2_w1),min(y3_w1)]);
maxY = max([max(y1_w1),max(y2_w1),max(y3_w1)]);
y_w1 = minY:0.05:maxY;

% for visualization lets compute the probability
% density function of the
% classes after projection
% the first class
y1_w1_Mu = mean(y1_w1);
y1_w1_sigma = std(y1_w1);
y1_w1_pdf = mvnpdf(y_w1',y1_w1_Mu,y1_w1_sigma);

% the second class
y2_w1_Mu = mean(y2_w1);
y2_w1_sigma = std(y2_w1);
y2_w1_pdf = mvnpdf(y_w1',y2_w1_Mu,y2_w1_sigma);

% the third class
y3_w1_Mu = mean(y3_w1);
y3_w1_sigma = std(y3_w1);
y3_w1_pdf = mvnpdf(y_w1',y3_w1_Mu,y3_w1_sigma);
```



Along second projection vector $y = w_2^T x$

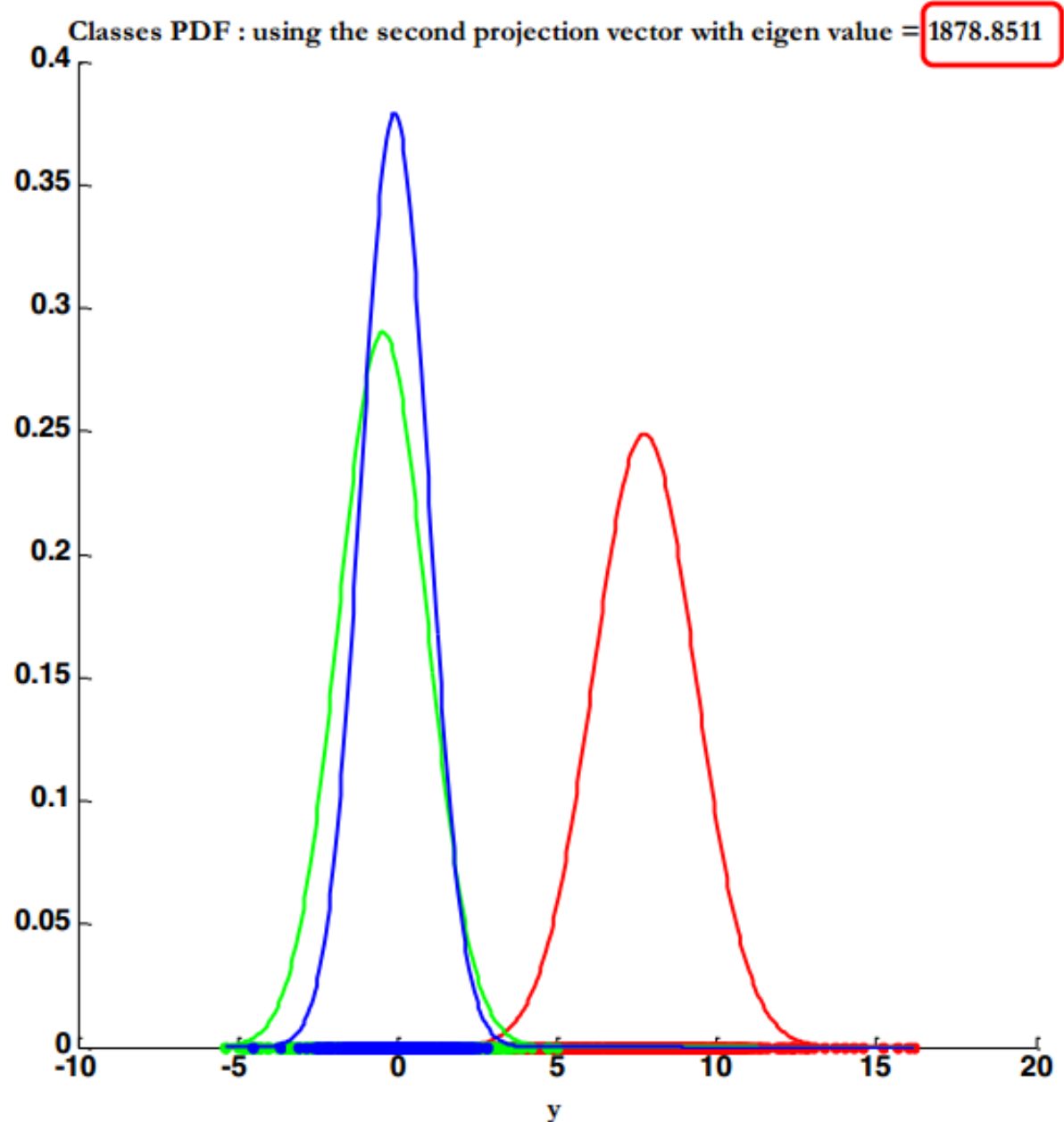
```
% project data samples along the projections axes
% the second projection vector
y1_w2 = W2'*X1;
y2_w2 = W2'*X2;
y3_w2 = W2'*X3;

% projection limits
minY = min([min(y1_w2),min(y2_w2),min(y3_w2)]);
maxY = max([max(y1_w2),max(y2_w2),max(y3_w2)]);
y_w2 = minY:0.05:maxY;

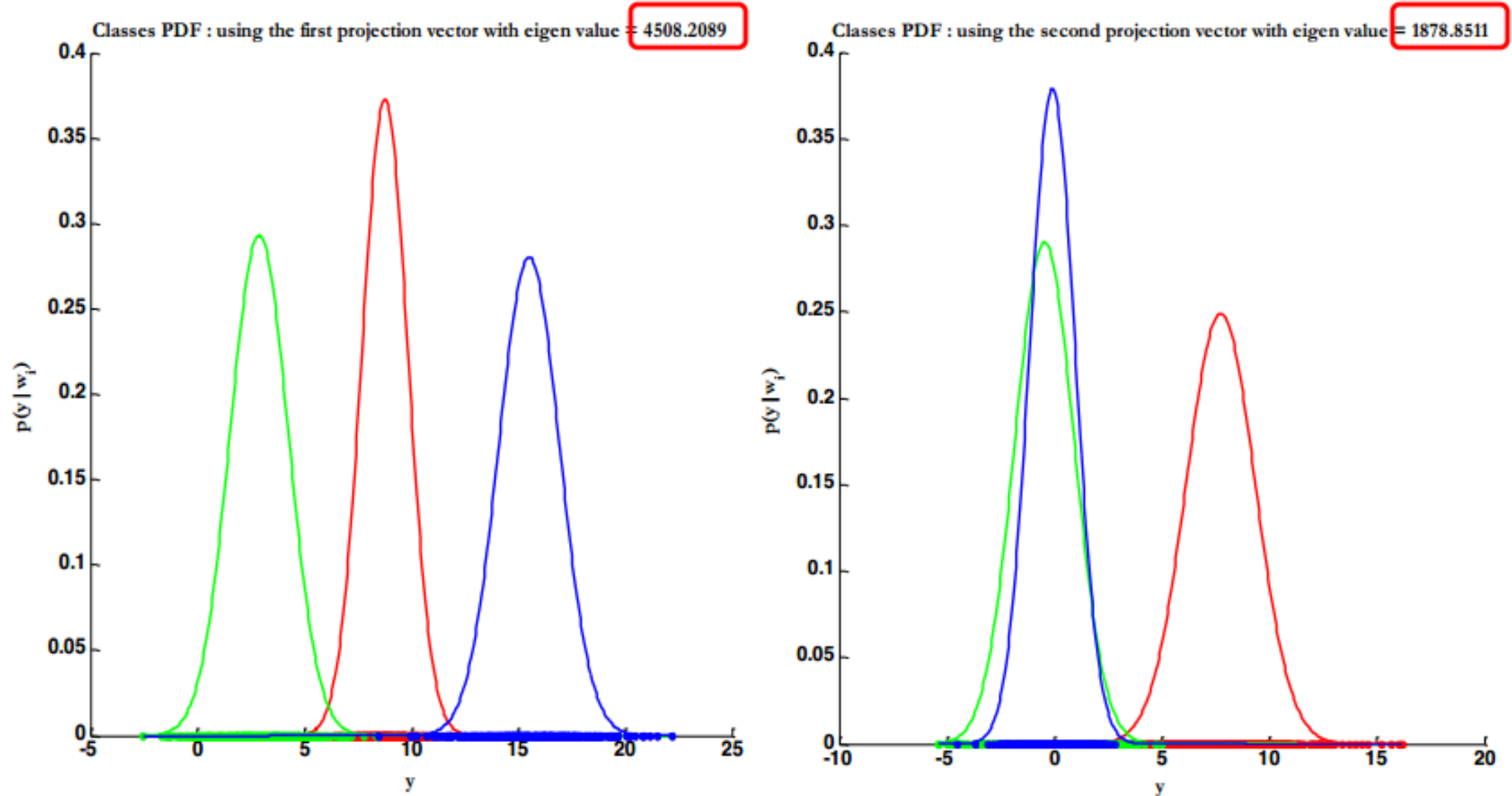
% for visualization lets compute the probability
% density function of the
% classes after projection
% the first class
y1_w2_Mu = mean(y1_w2);
y1_w2_sigma = std(y1_w2);
y1_w2_pdf = mvnpdf(y_w2',y1_w2_Mu,y1_w2_sigma);

% the second class
y2_w2_Mu = mean(y2_w2);
y2_w2_sigma = std(y2_w2);
y2_w2_pdf = mvnpdf(y_w2',y2_w2_Mu,y2_w2_sigma);

% the third class
y3_w2_Mu = mean(y3_w2);
y3_w2_sigma = std(y3_w2);
y3_w2_pdf = mvnpdf(y_w2',y3_w2_Mu,y3_w2_sigma);
```



Apparently, the projection vector that has the highest eigenvalue provides higher discrimination power between classes.



Classification with LDA

- First, select the number of feature dimension for the low-dimensional feature space. (In most cases, LDA is used for **dimension reduction**.)
- Nearest Neighbor or other classifiers.
- In practice, mostly PCA shall be performed before LDA to avoid the singularity issue.

Summary

- Linear Discriminant Analysis—Two Classes
 - Minimize within-class scatter
 - Maximize between-class scatter
 - The eigenvector of the **largest** eigenvalue of $\mathbf{S}_w^{-1}\mathbf{S}_b$ (as $-\boldsymbol{\theta}^{*T}\mathbf{S}_b\boldsymbol{\theta}^* = -\lambda\boldsymbol{\theta}^{*T}\mathbf{S}_w\boldsymbol{\theta}^* = -\lambda$)
 - Or $\boldsymbol{\theta}^* = \mathbf{S}_w^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$
- Linear Discriminant Analysis— C Classes
 - **Dimension reduction.** $\Theta \in \mathbb{R}^{d \times p} : \mathbf{X} \rightarrow \mathbf{Y}$ ($p \ll d$). In fact, $p \leq C + 1$.
 - Columns of Θ^* are eigenvectors of $\mathbf{S}_w^{-1}\mathbf{S}_b$ corresponding to the p largest eigenvalues.

Statistical Facts

Between-class scatter:

$$\mathbf{S}_b = \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T = \frac{1}{2N} \sum_{i,j=1}^C n_i n_j (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T$$

$$\begin{aligned} \frac{1}{2N} \sum_{i,j=1}^C n_i n_j (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)(\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)^T &= \frac{1}{2N} \sum_{i,j=1}^C n_i n_j [(\boldsymbol{\mu}_i - \boldsymbol{\mu}) + (\boldsymbol{\mu} - \boldsymbol{\mu}_j)][(\boldsymbol{\mu}_i - \boldsymbol{\mu}) + (\boldsymbol{\mu} - \boldsymbol{\mu}_j)]^T \\ &= \frac{1}{2N} \sum_{i,j=1}^C n_i n_j [(\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + (\boldsymbol{\mu} - \boldsymbol{\mu}_j)(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu} - \boldsymbol{\mu}_j)^T + (\boldsymbol{\mu} - \boldsymbol{\mu}_j)(\boldsymbol{\mu} - \boldsymbol{\mu}_j)^T] \\ &= \frac{1}{2N} \sum_{i,j=1}^C n_i n_j [(\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + (\boldsymbol{\mu} - \boldsymbol{\mu}_j)(\boldsymbol{\mu} - \boldsymbol{\mu}_j)^T] \\ &= \frac{1}{2} \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T + \frac{1}{2} \sum_{j=1}^C n_j (\boldsymbol{\mu} - \boldsymbol{\mu}_j)(\boldsymbol{\mu} - \boldsymbol{\mu}_j)^T \\ &= \sum_{i=1}^C n_i (\boldsymbol{\mu}_i - \boldsymbol{\mu})(\boldsymbol{\mu}_i - \boldsymbol{\mu})^T = \mathbf{S}_b \end{aligned}$$