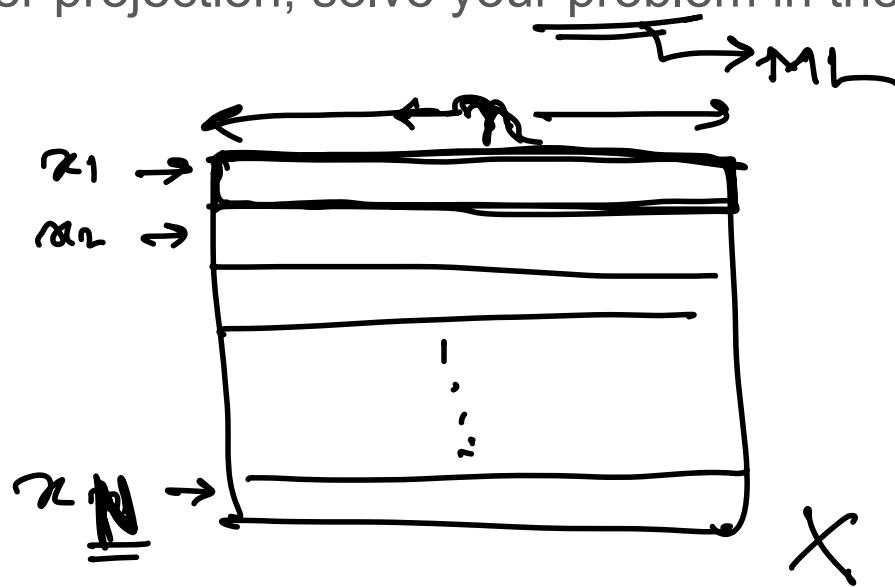
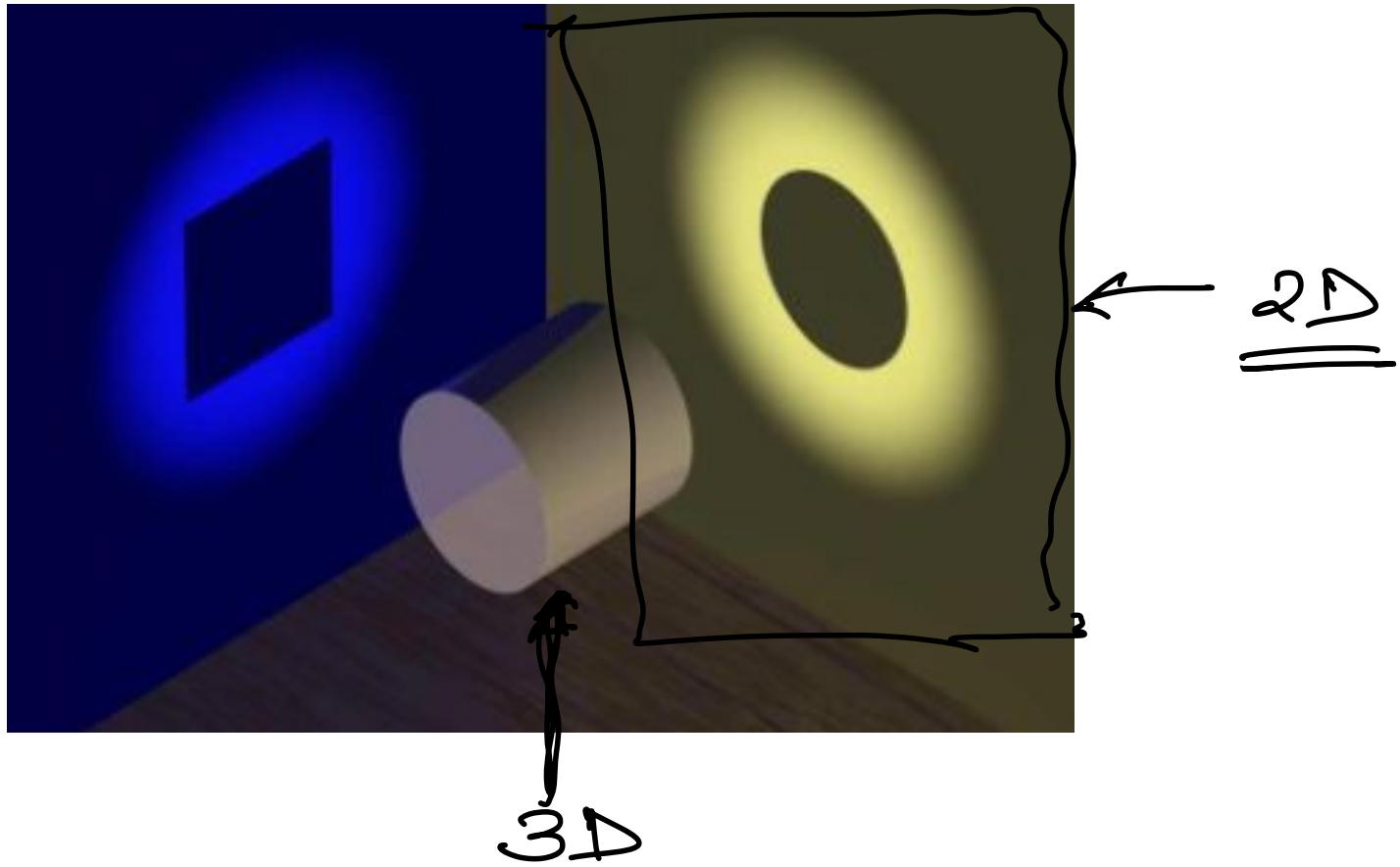


Dimensionality Reduction

- The main idea – Reduce the dimensionality of the space
- Project the n dimensional points to a k dimensional space so that
 - $k \ll n$
 - Some important properties of the data are preserved
- After projection, solve your problem in the low dimensional space

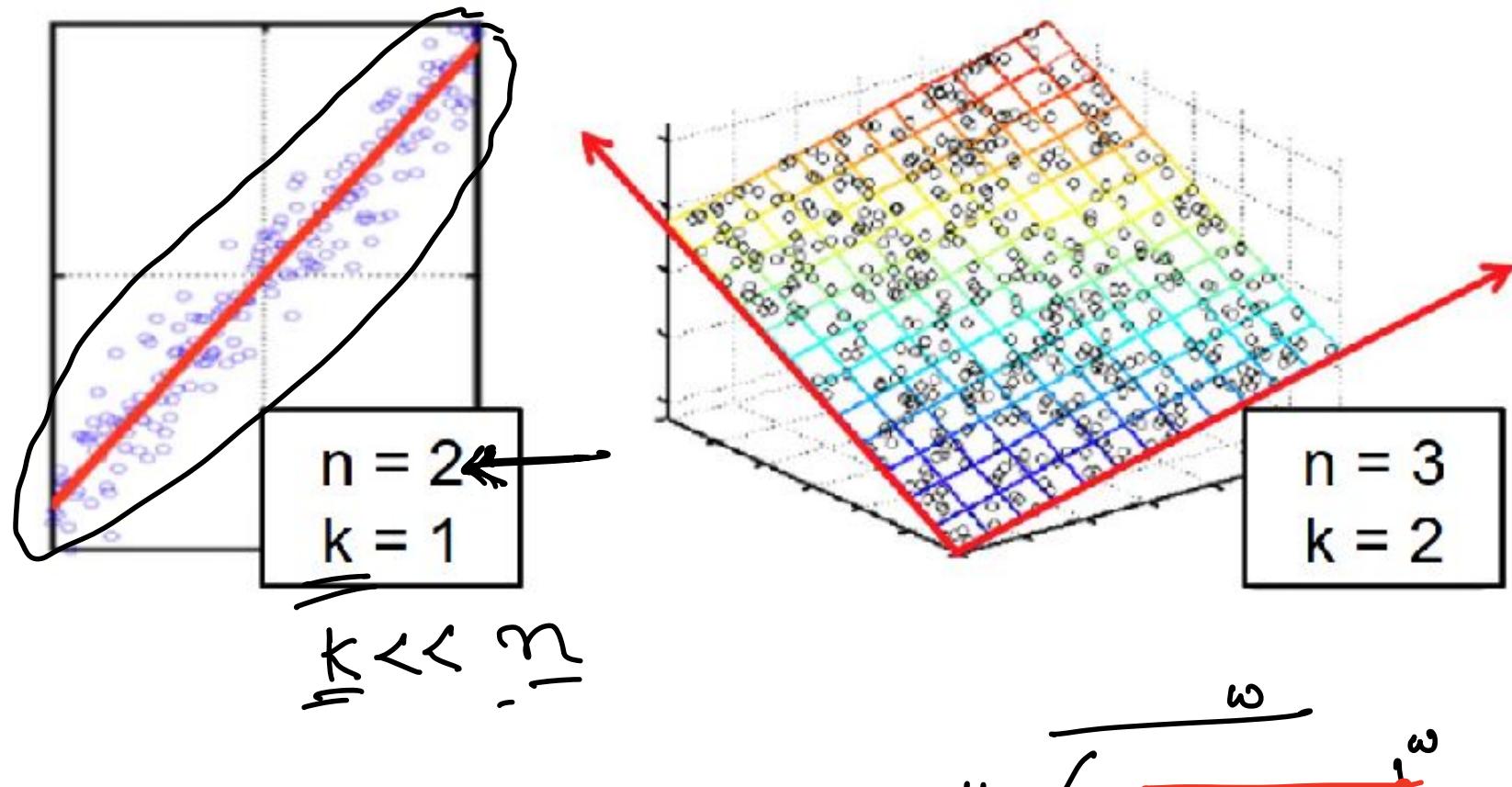




Dimensionality Reduction (benefits)

- Easier learning – fewer parameters
 - Less computation time for training/inference
 - Less space needed to store data
- Visualization – show high dimensional data in 2-d
- Discover “intrinsic dimensionality” of data
 - noise removal
 - high dimensional data that is actually lower dimensional

High dimensional data that is actually lower dimensional

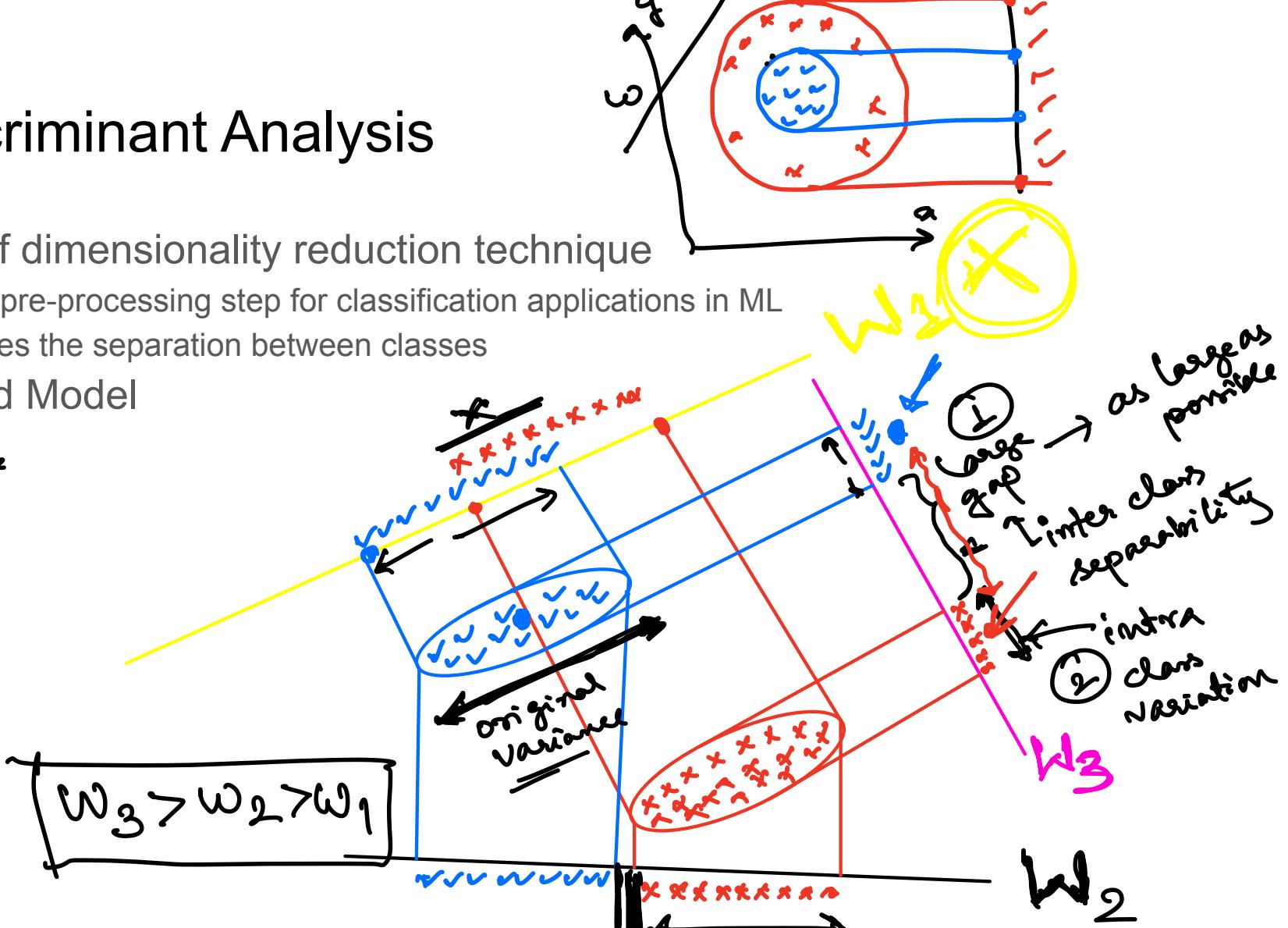


Linear Discriminant Analysis

- One type of dimensionality reduction technique
 - used as pre-processing step for classification applications in ML
 - maximizes the separation between classes
- Supervised Model

(Q) Original variance
v/s
Projected variance

(Q) When is LDA not
a good option





$w_2 > w_1$

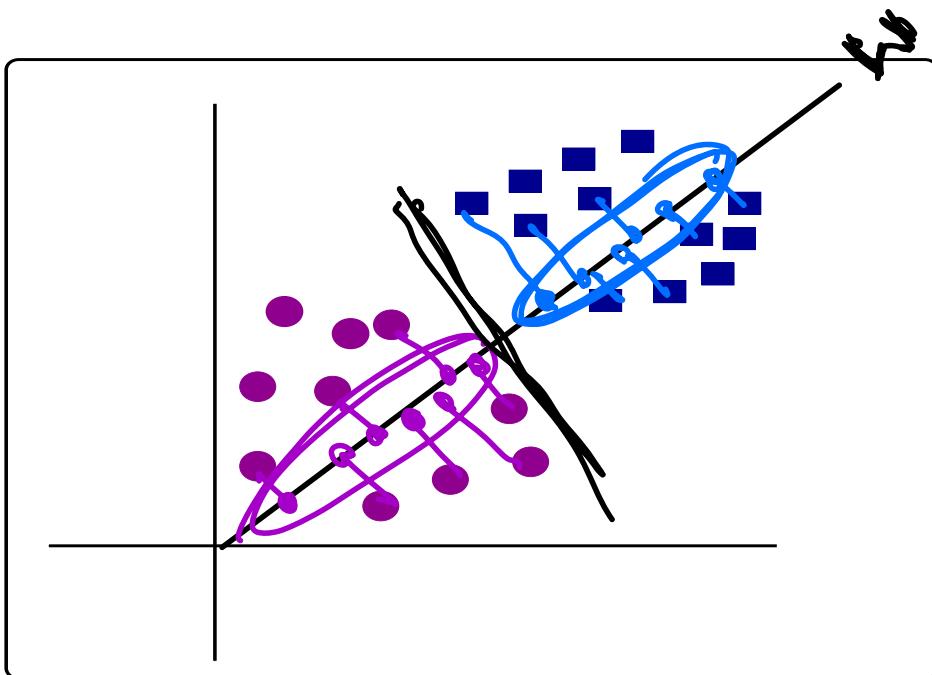
Fisher Linear Discriminant

- A linear discriminant function based classifier is:

Decide $\underline{\underline{X}} \in C-1$ if $\underline{\underline{W}}^T \underline{\underline{X}} + \underline{\underline{w}_0} > 0$

- Hence One can think of the best W as the direction along which the two classes are well separated.
- We project the data along the direction $\underline{\underline{W}}$. Separation between points of different classes in the projected data is a good way to rate how good is W .
- Such a method is called Fisher Linear Discriminant.

- Consider the following 2-class example



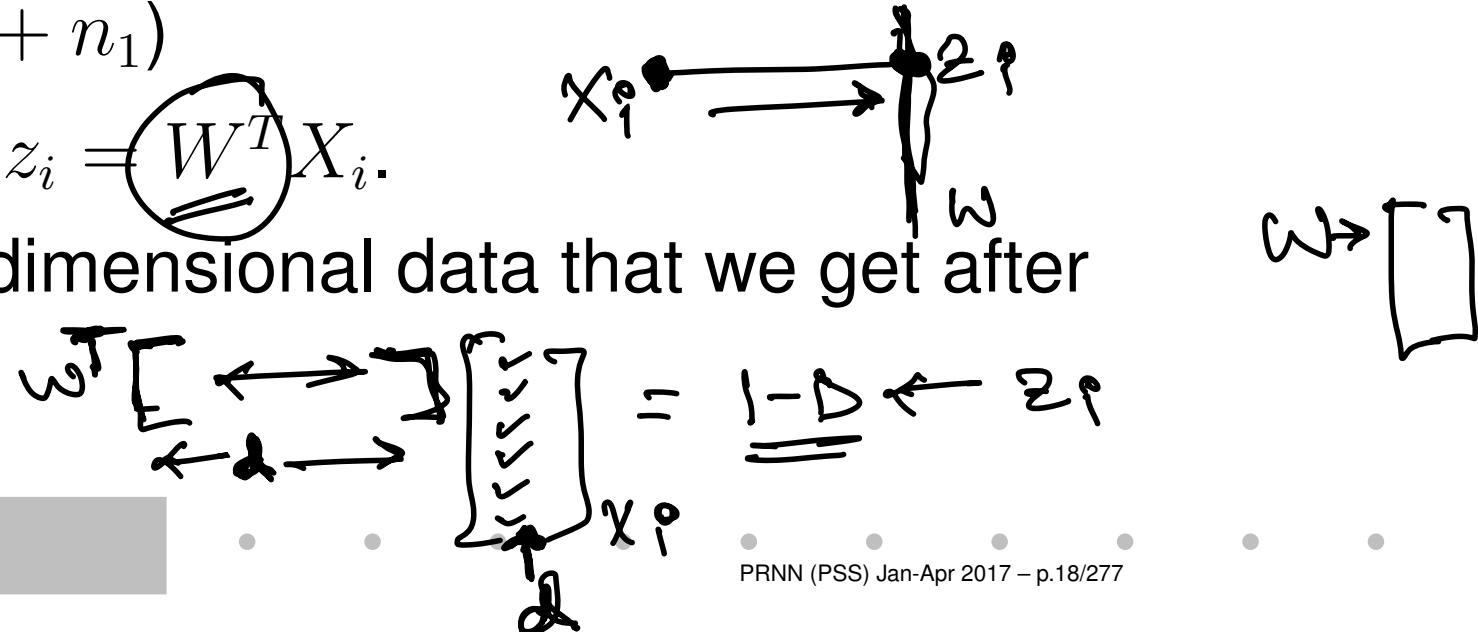
•
•

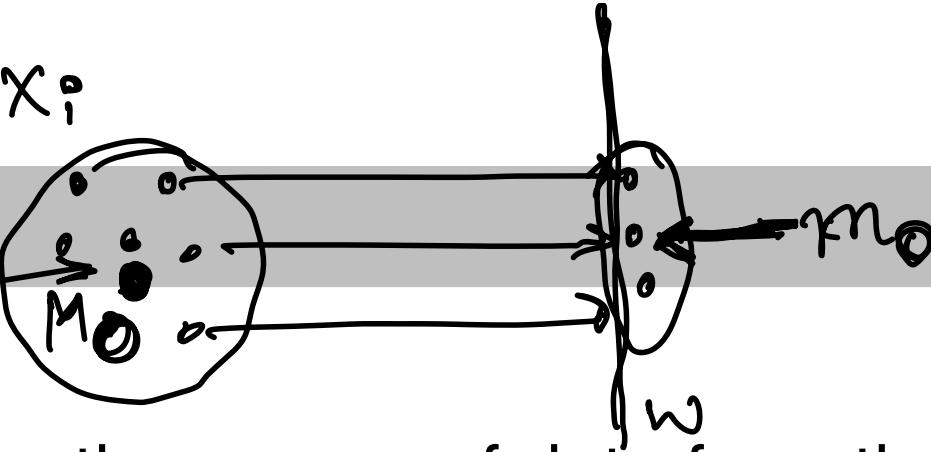
Fisher Linear Discriminant

- The idea is to find a direction W such that the training data of the two classes are well-separated if projected onto this direction.
 - We need some figure of merit for each W to characterize how well the W results in such a separation.
 - We consider the 2-class case.
-

+training data

- Let $\{(X_i, y_i), i = 1, \dots, n\}$ be the data.
- Let $y_i \in \{0, 1\}$.
- Let C_0 and C_1 denote the two classes. Thus, if $y_i = 0$ then $X_i \in C_0$ and if $y_i = 1$ then $X_i \in C_1$.
- Let n_0 and n_1 denote the number of examples of each class. ($n = n_0 + n_1$)
- For any W , let $z_i = W^T X_i$.
- z_i are the one dimensional data that we get after projection.





- Let M_0 and M_1 be the means of data from the two classes:

$$M_0 = \frac{1}{n_0} \sum_{X_i \in C_0} X_i; \quad M_1 = \frac{1}{n_1} \sum_{X_i \in C_1} X_i$$

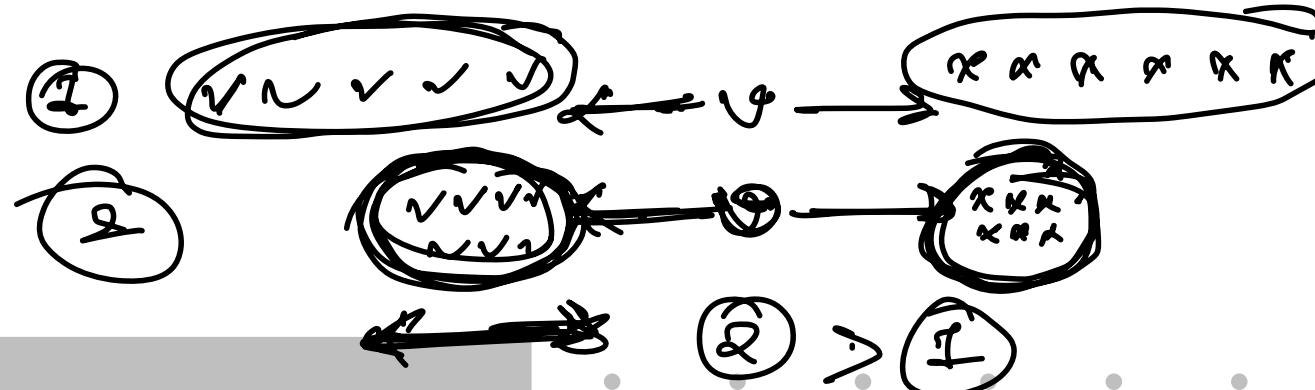
The corresponding means of the projected data would be

$$m_0 = W^T M_0 \quad \text{and} \quad m_1 = W^T M_1$$

$$\frac{1}{n_0} \sum_{x_i \in C_0} \overbrace{W^T x_i}^{=} = \frac{1}{n_0} \sum_{x_i \in C_0} x_i$$

$$= \overbrace{W^T}^{\cdot} \overbrace{M_0}^{\cdot}$$

- The difference $(m_0 - m_1)$ gives us an idea of the separation between samples of the two classes after projecting the data onto the direction W .
- Hence, we may want a W that maximizes $(m_0 - m_1)^2$.
- However, we have to make this scale independent.
- Also, the distance between means should be viewed relative to the variances. 



- Define

$$s_0^2 = \sum_{X_i \in C_0} (W^T X_i - m_0)^2; \quad s_1^2 = \sum_{X_i \in C_1} (W^T X_i - m_1)^2$$

↑

These give us the variances (upto a factor) of the two classes in the projected data.

- We want large separation between m_0 and m_1 relative to the variances.

- Hence we can take our objective to be to maximize

$$J(W) = \frac{(m_1 - m_0)^2}{s_0^2 + s_1^2}$$

(m₁ - m₀)² ← }
 s₀² + s₁² ← }

- We now rewrite J into a more convenient form.
- We have

$$\begin{aligned}
 \underline{(m_1 - m_0)^2} &= \underline{(W^T M_1 - W^T M_0)^2} \\
 &= W^T \underline{(M_1 - M_0)(M_1 - M_0)^T} W \\
 &\quad - \rightarrow S_B
 \end{aligned}$$

$\stackrel{\longleftrightarrow}{=} \stackrel{\longrightarrow}{[W^T(M_1 - M_0)]^2}$

$$x^2 = x^T x$$

$$\varphi^T (M_1 - M_0) (M_1 - M_0)^T \varphi = \underline{\underline{||\varphi^T (M_1 - M_0)||^2}} \geq 0$$

$$\varphi^T S_B \varphi \geq 0 \quad \forall \varphi \in \mathbb{R}^d$$

- Thus we have $(m_1 - m_0)^2 = \underline{\underline{W^T S_B W}}$ where

Symmetric \rightarrow P.S.d.

$$\underline{\underline{\text{Rank 1}}} \rightarrow \underline{\underline{S_B}} = \underline{\underline{(M_1 - M_0)(M_1 - M_0)^T}}.$$

- Here, S_B is a $d \times d$ matrix (note that $X_i \in \mathbb{R}^d$).

- It is called between class scatter matrix.

- We can similarly write s_0^2 and s_1^2 also as quadratic forms.

$$x_i \rightarrow \mathbb{R}^d$$

$$M_1, M_0 \rightarrow \mathbb{R}^d$$

$$\begin{bmatrix} 1 \\ \vdots \\ d \end{bmatrix} \left[\begin{array}{c} \nearrow d \\ \searrow \end{array} \right] \downarrow$$

$$\begin{bmatrix} 1 \\ 2 \\ 3 \end{bmatrix} \begin{bmatrix} 1 & 2 & 3 \end{bmatrix} = \begin{bmatrix} 1 & 2 & 3 \\ 2 & 4 & 6 \\ 3 & 6 & 9 \end{bmatrix}$$

We have

$$\begin{aligned}s_0^2 &= \sum_{X_i \in C_0} (W^T \underline{\underline{X}}_i - \underline{\underline{\underline{W}}^T M}_0)^2 \\&= \sum_{X_i \in C_0} [W^T(X_i - M_0)]^2 \\&= \sum_{X_i \in C_0} W^T(X_i - M_0)(X_i - M_0)^T W \\&= \underline{\underline{W}}^T \left[\underbrace{\sum_{X_i \in C_0} (X_i - M_0)(X_i - M_0)^T}_{\text{---}} \right] \underline{\underline{W}}\end{aligned}$$

- Similarly, we get

$$\underline{s_1^2} = W^T \left[\sum_{X_i \in C_1} (X_i - M_1)(X_i - M_1)^T \right] W$$

- Thus we can write $s_0^2 + s_1^2 = W^T S_w W$, where

$$S_w = \sum_{X_i \in C_0} (X_i - M_0)(X_i - M_0)^T + \sum_{X_i \in C_1} (X_i - M_1)(X_i - M_1)^T$$

- S_w is also $d \times d$ matrix and is called *within class scatter matrix*.

• $\|S_w\| \neq \mathcal{D}$

$$w \rightarrow \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix}_d$$

- Hence we can now write J as

$$\text{Scalar} \rightarrow J(W) = \frac{W^T S_B W}{W^T S_w W} \xrightarrow{\cancel{d \times d}} \frac{\text{Scalar}}{(m_0 - m_1)^2 / (S_0^2 + S_1^2)}$$

- We want to find a \underline{W} that maximizes $\underline{J(W)}$.
- Note that $J(W)$ is not affected by scaling of W . \circled{w}
- Given the data we can calculate the S_B and S_w .
- Maximizing ratio of quadratic forms is a standard optimization problem.

- We need to maximize

$$J(W) = \frac{W^T S_B W}{W^T S_w W}$$

~~$\frac{\partial (w^T S_B w)}{\partial w}$~~
 $+ w^T S_B w \frac{\partial (\frac{1}{w^T S_B w})}{\partial w}$

- Differentiating w.r.t. W and equating to zero, we get

$$\cancel{(w^T S_w w)} \cancel{\left[\frac{2S_B W}{W^T S_w W} \right]} - \left[\frac{W^T S_B W}{(W^T S_w W)^2} \cancel{2S_w W} \right] = 0 \quad (w^T S_w w)$$

- Implies, $S_B W$ is in the same direction as $S_w W$.

$$w^T S_w w \neq 0$$

$$S_B W = \boxed{\frac{w^T S_B w}{w^T S_w w}} S_w W = 0$$

PBNN / PSS | Jan-Apr 2017 – p.49/277

$$S_B \omega - J(\omega) S_\omega \omega = 0$$

Scalar

- Thus, any maximizer of $J(W)$ has to satisfy

$$\boxed{S_w W} = \cancel{\lambda} S_B W \iff S_w^{-1}$$

exists

for some constant λ .

- This is known as the generalized eigen value problem.
- There are standard methods to solve this problem using, e.g., LU decomposition.
- By solving the generalized eigen value problem we can find the best direction W .

- Often, the real symmetric matrix S_w would be invertible.
- Recall that

$$S_w = \sum_{X_i \in C_0} (X_i - M_0)(X_i - M_0)^T + \sum_{X_i \in C_1} (X_i - M_1)(X_i - M_1)^T$$

- This is a sum of large number of rank 1 matrices.

- If S_w is invertible, then we can write

$$\cancel{S_w^{-1} S_w} W = \cancel{S_w^{-1} S_B} W$$

$$\underline{\underline{m_0}} = \underline{\underline{w^T M_0}}$$

- We have

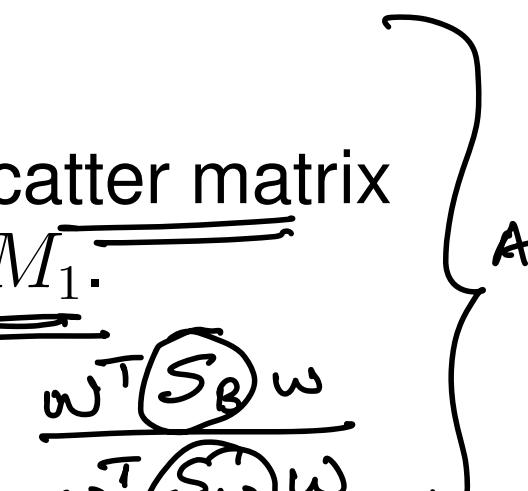
$$S_B W = (M_1 - M_0) \cancel{(M_1 - M_0)^T W} = k(M_1 - M_0)$$

where k is some constant. (note $k = (m_1 - m_0)$) \rightarrow scalar

- Now we get (since scale factor in W is not relevant)

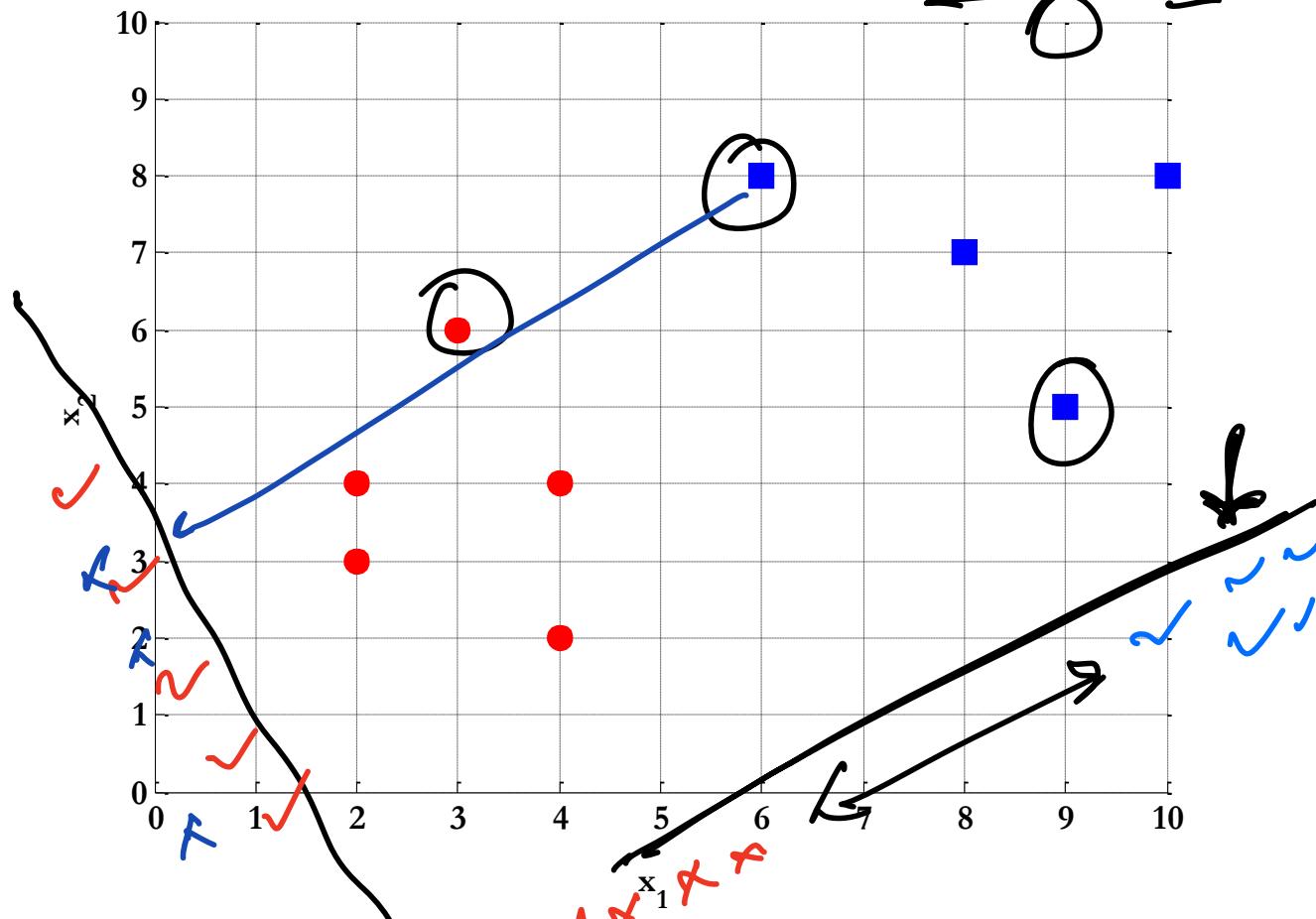
$$\underline{\underline{W}} = S_w^{-1} (\cancel{M_1 - M_0}) \cancel{\lambda} \cancel{k}$$

Obtaining Fisher Linear Discriminant

- We can sum-up the process as follows.
 - Given the training data, we first form the scatter matrix S_w and also calculate the means M_0 and M_1 .
 - If S_w is invertible, we calculate W by
$$W = S_w^{-1}(M_1 - M_0)$$
 - Even if S_w is not invertible, there are techniques to find the maximizer of $J(W)$ by solving the generalized eigen value problem.
 - Thus we can find the best direction W .
- Alg
- 

LDA ... Two Classes - Example

- Compute the Linear Discriminant projection for the following two-dimensional dataset.
 - Samples for class ω_1 : $\mathbf{X}_1 = (x_1, x_2) = \{(4,2), (2,4), (2,3), (3,6), (4,4)\}$ \leftarrow Red
 - Sample for class ω_2 : $\mathbf{X}_2 = (x_1, x_2) = \{(9,10), (6,8), (9,5), (8,7), (10,8)\}$ \leftarrow Blue



```
% samples for class 1
X1 = [4,2;
       2,4;
       2,3;
       3,6;
       4,4];

% samples for class 2
X2 = [9,10;
       6,8;
       9,5;
       8,7;
       10,8];
```

LDA ... Two Classes - Example

- The classes mean are :

$$\underline{\mu}_1 = \frac{1}{N_1} \sum_{x \in \omega_1} x = \frac{1}{5} \left[\begin{pmatrix} 4 \\ 2 \end{pmatrix} + \begin{pmatrix} 2 \\ 4 \end{pmatrix} + \begin{pmatrix} 2 \\ 3 \end{pmatrix} + \begin{pmatrix} 3 \\ 6 \end{pmatrix} + \begin{pmatrix} 4 \\ 4 \end{pmatrix} \right] = \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \leftarrow$$

$$\underline{\mu}_2 = \frac{1}{N_2} \sum_{x \in \omega_2} x = \frac{1}{5} \left[\begin{pmatrix} 9 \\ 10 \end{pmatrix} + \begin{pmatrix} 6 \\ 8 \end{pmatrix} + \begin{pmatrix} 9 \\ 5 \end{pmatrix} + \begin{pmatrix} 8 \\ 7 \end{pmatrix} + \begin{pmatrix} 10 \\ 8 \end{pmatrix} \right] = \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \leftarrow$$

```
% class means  
Mu1 = mean(X1)';  
Mu2 = mean(X2)';
```

LDA ... Two Classes - Example

- Covariance matrix of the first class:

$$S_1 = \sum_{x \in \omega_1} (x - \underline{\mu}_1)(x - \underline{\mu}_1)^T = \left[\begin{pmatrix} 4 \\ 2 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 2 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2$$
$$+ \left[\begin{pmatrix} 2 \\ 3 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 3 \\ 6 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2 + \left[\begin{pmatrix} 4 \\ 4 \end{pmatrix} - \begin{pmatrix} 3 \\ 3.8 \end{pmatrix} \right]^2$$

d=2

$$= \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} \leftarrow \text{intra class scatter matrix}$$

```
% covariance matrix of the first class  
S1 = cov(X1);
```

LDA ... Two Classes - Example

- Covariance matrix of the second class:

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

```
% covariance matrix of the first class  
S2 = cov(X2);
```

LDA ... Two Classes - Example

- Within-class scatter matrix:

$$\cancel{S_w = S_1 + S_2} = \begin{pmatrix} 1 & -0.25 \\ -0.25 & 2.2 \end{pmatrix} + \begin{pmatrix} 2.3 & -0.05 \\ -0.05 & 3.3 \end{pmatrix}$$
$$= \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix} \leftarrow \boxed{S_w^{-1}(M_1 - M_0)}$$

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

LDA ... Two Classes - Example

- Between-class scatter matrix:

$$S_B = (\underline{\mu_1} - \underline{\mu_2})(\underline{\mu_1} - \underline{\mu_2})^T$$

$$\boxed{S_w^{-1} S_B} = 1$$

$\rightarrow \text{rank}(AB) \leq \min \{ \text{rank}(A), \text{rank}(B) \}$

$$= \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right] \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]^T$$

$$= \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix} \begin{pmatrix} -5.4 & -3.8 \end{pmatrix}$$

Inter class Scatter

$$= \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix}$$

```
% between-class scatter matrix
SB = (Mu1-Mu2) * (Mu1-Mu2)';
```

LDA ... Two Classes - Example

- The LDA projection is then obtained as the solution of the generalized eigen value problem

$$\begin{aligned}
 S_W^{-1} S_B w = \lambda w &\rightarrow S_w^{-1} S_B w - \lambda w = 0 \\
 \Rightarrow |S_W^{-1} S_B - \lambda I| = 0 & \quad (S_w^{-1} S_B - \lambda I) \underline{w} = 0 \\
 \Rightarrow \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix}^{-1} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0 \\
 \Rightarrow \begin{pmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{pmatrix} \begin{pmatrix} 29.16 & 20.52 \\ 20.52 & 14.44 \end{pmatrix} - \lambda \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix} = 0 \\
 \Rightarrow \begin{pmatrix} 9.2213 - \lambda & 6.489 \\ 4.2339 & 2.9794 - \lambda \end{pmatrix} = 0 \\
 \Rightarrow (9.2213 - \lambda)(2.9794 - \lambda) - 6.489 \times 4.2339 = 0 \\
 \Rightarrow \lambda^2 - 12.2007\lambda = 0 \Rightarrow \lambda(\lambda - 12.2007) = 0 \\
 \Rightarrow \lambda_1 = 0, \lambda_2 = 12.2007
 \end{aligned}$$

LDA ... Two Classes - Example

- Hence

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

and

$$A = S_w^{-1} S_B$$

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

Thus;

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

and

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

```
% computing the LDA projection
invSw = inv(Sw);

invSw_by_SB = invSw * SB;

% getting the projection vector
[V,D] = eig(invSw_by_SB)

% the projection vector
W = V(:,1);
```

- The optimal projection is the one that given maximum $\lambda = J(w)$

LDA ... Two Classes - Example

Or directly;

S_w^{-1} is invertible

$$w^* = S_W^{-1} \left(\underline{\mu_1} - \underline{\mu_2} \right) = \begin{pmatrix} 3.3 & -0.3 \\ -0.3 & 5.5 \end{pmatrix}^{-1} \left[\begin{pmatrix} 3 \\ 3.8 \end{pmatrix} - \begin{pmatrix} 8.4 \\ 7.6 \end{pmatrix} \right]$$

$\cancel{S_w^{-1}}$ easier

$$\begin{bmatrix} a & b \\ c & d \end{bmatrix}^{-1} = \begin{bmatrix} d & -b \\ -c & a \end{bmatrix}$$

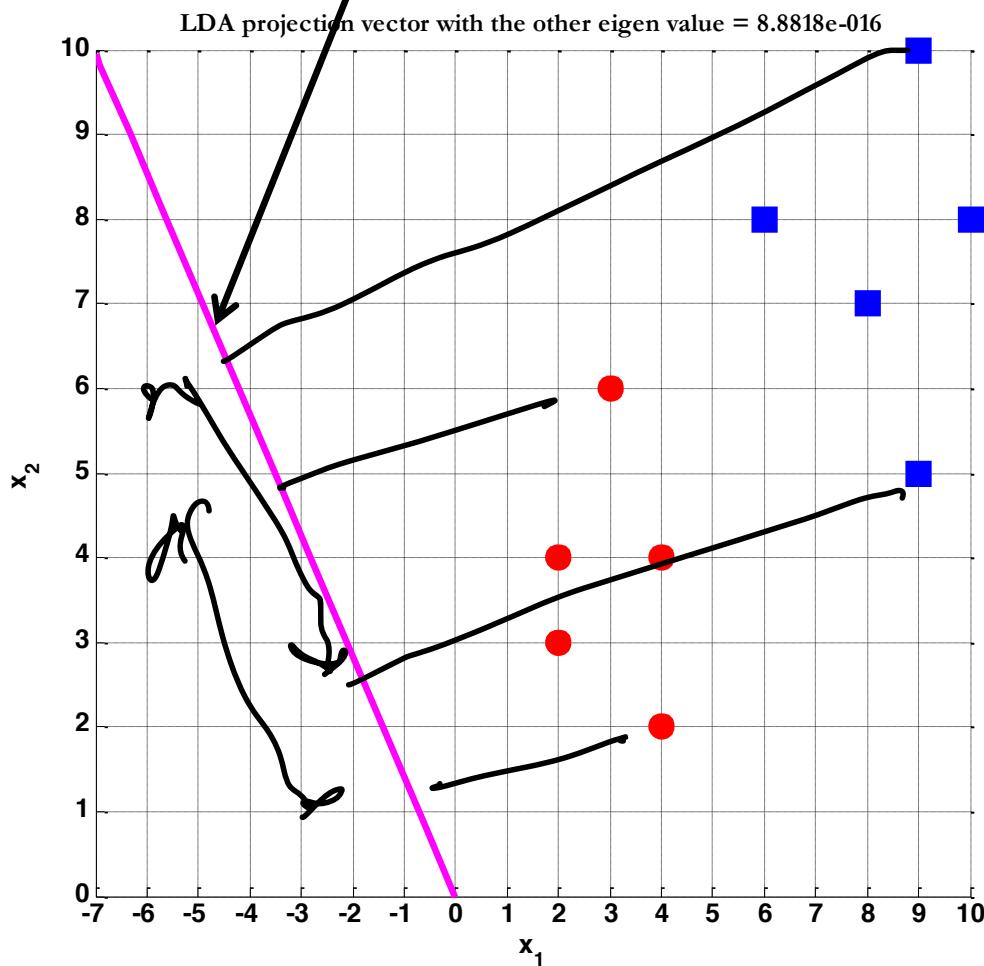
$$= \begin{pmatrix} 0.3045 & 0.0166 \\ 0.0166 & 0.1827 \end{pmatrix} \begin{pmatrix} -5.4 \\ -3.8 \end{pmatrix}$$
$$= \begin{pmatrix} 0.9088 \\ 0.4173 \end{pmatrix}$$

$\cancel{w^*}$

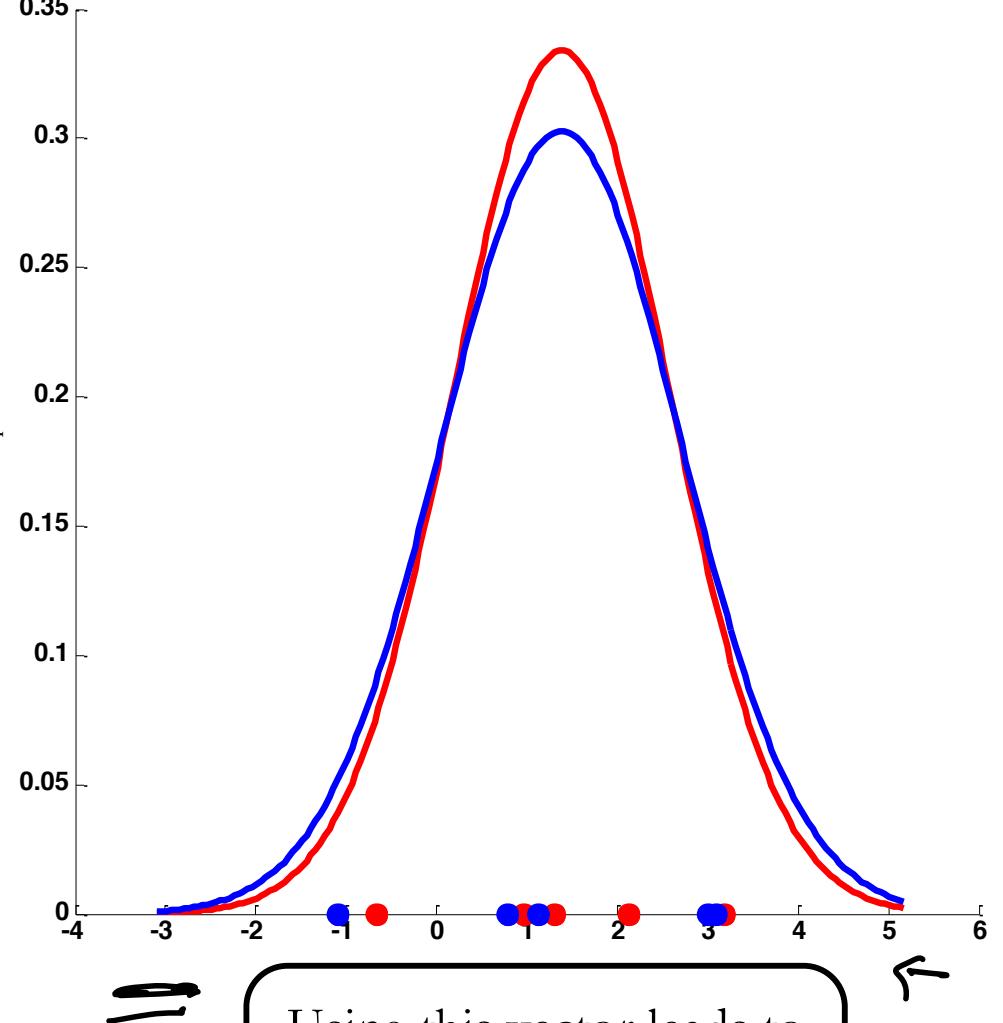
LDA - Projection

The projection vector corresponding to the **smallest** eigen value

$$\lambda = 0$$



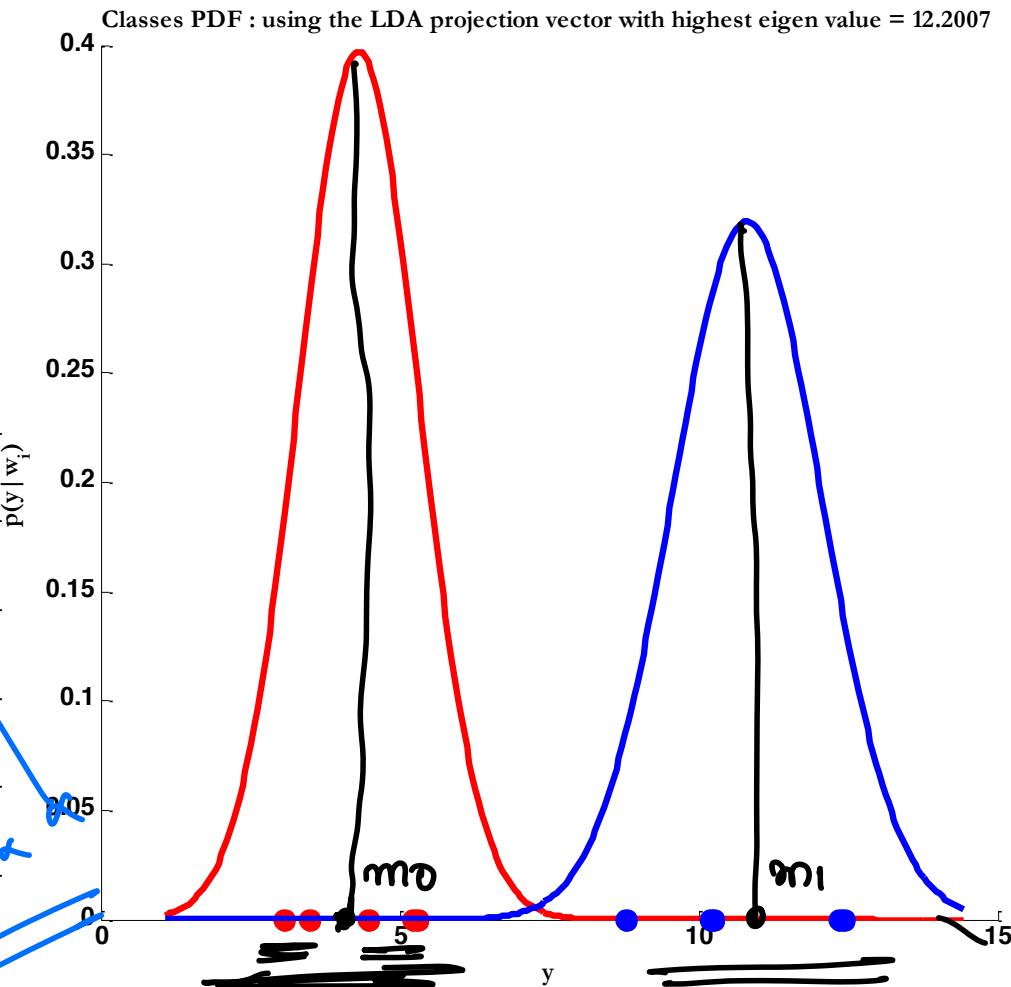
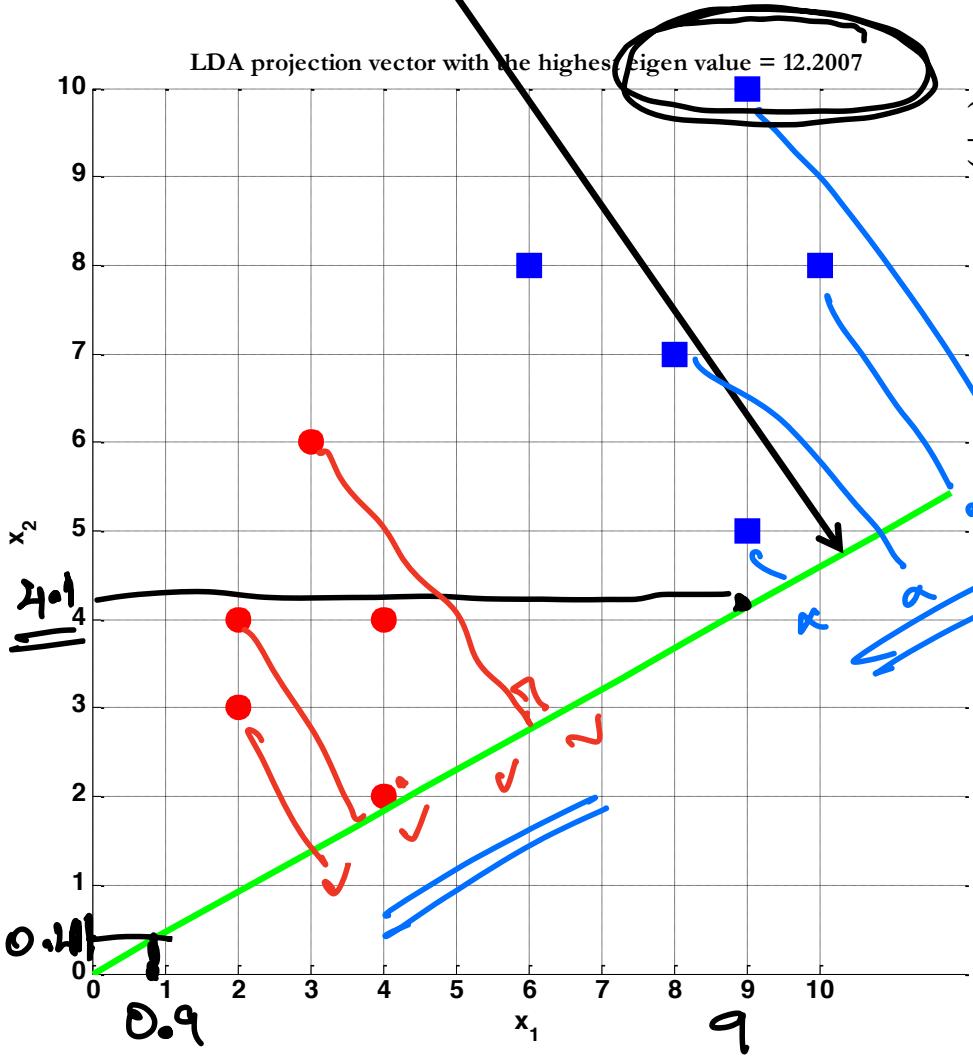
Classes PDF : using the LDA projection vector with the other eigen value = 8.8818e-016



Using this vector leads to
bad separability
between the two classes

LDA - Projection

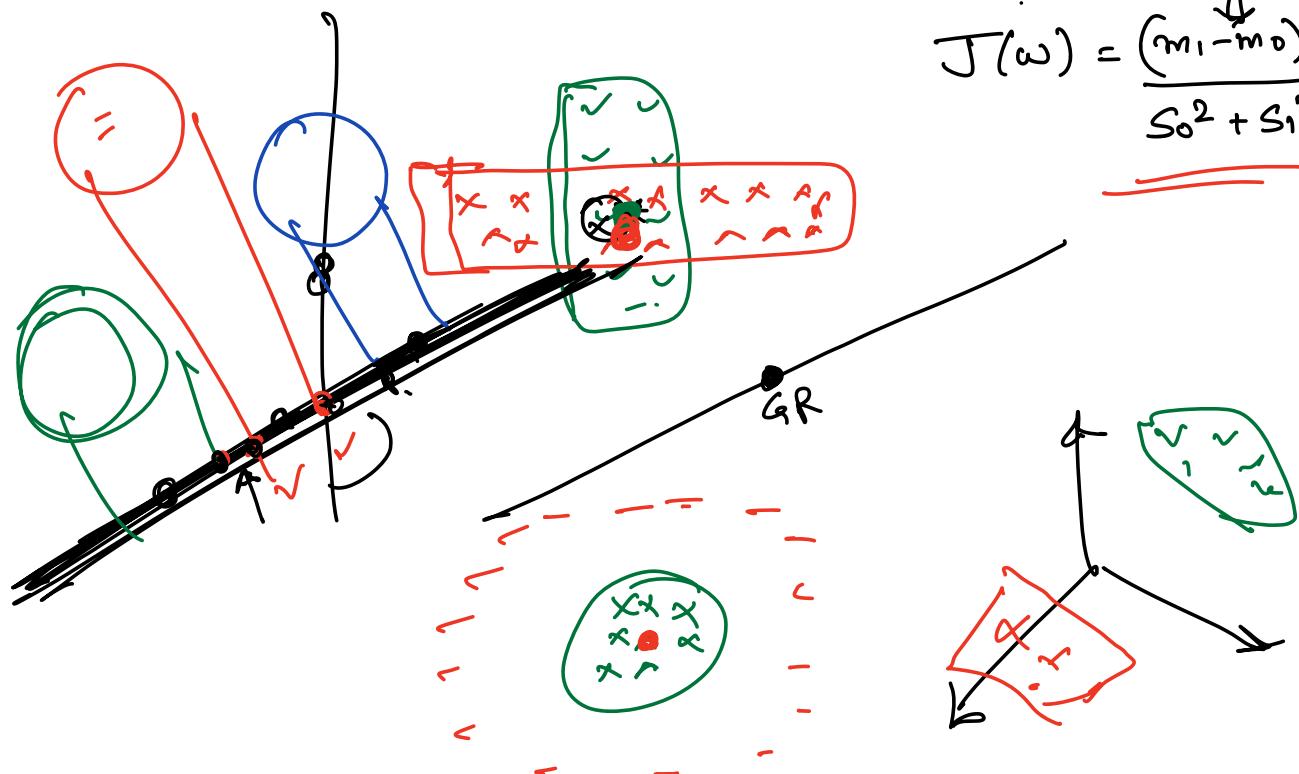
The projection vector corresponding to the **highest** eigen value



Using this vector leads to
good separability
between the two classes



$$J(\omega) = \frac{(m_1 - m_0)^2}{S_0^2 + S_1^2}$$



LDA ... C-Classes

- Now, we have C -classes instead of just two.
- We are now seeking $(C-1)$ projections $[y_1, y_2, \dots, y_{C-1}]$ by means of $(C-1)$ projection vectors w_i .
- w_i can be arranged by *columns* into a projection matrix $W = [w_1 | w_2 | \dots | w_{C-1}]$ such that:

$$\underbrace{y_i}_{\substack{\nearrow \\ \text{projection}}} = \underbrace{w_i^T}_{=} \underbrace{x}_{=} \Rightarrow \underbrace{y}_{=} = \underbrace{W^T}_{=} \underbrace{x}_{=}$$

where $\underbrace{x}_{m \times 1} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$, $\underbrace{y}_{C-1 \times 1} = \begin{bmatrix} y_1 \\ \vdots \\ y_{C-1} \end{bmatrix}$

and $W_{m \times C-1} = [w_1 | w_2 | \dots | w_{C-1}]$

LDA ... C-Classes

- If we have n -feature vectors, we can stack them into one matrix as follows;

$$Y = W^T X$$

\equiv

feature matrix

where $X_{m \times n} = \begin{bmatrix} x_1^1 & x_1^2 & \dots & x_1^n \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ x_m^1 & x_m^2 & \dots & x_m^n \end{bmatrix}$, $Y_{C-1 \times n} = \begin{bmatrix} y_1^1 & y_1^2 & \dots & y_1^n \\ \vdots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ y_{C-1}^1 & y_{C-1}^2 & \dots & y_{C-1}^n \end{bmatrix}$

and $W_{m \times C-1} = [w_1 | w_2 | \dots | w_{C-1}]$

$\equiv \quad = \quad =$

$\xleftarrow{\text{C-1 dim}}$

LDA – C-Classes

- Recall the two classes case, the *within-class scatter* was computed as:

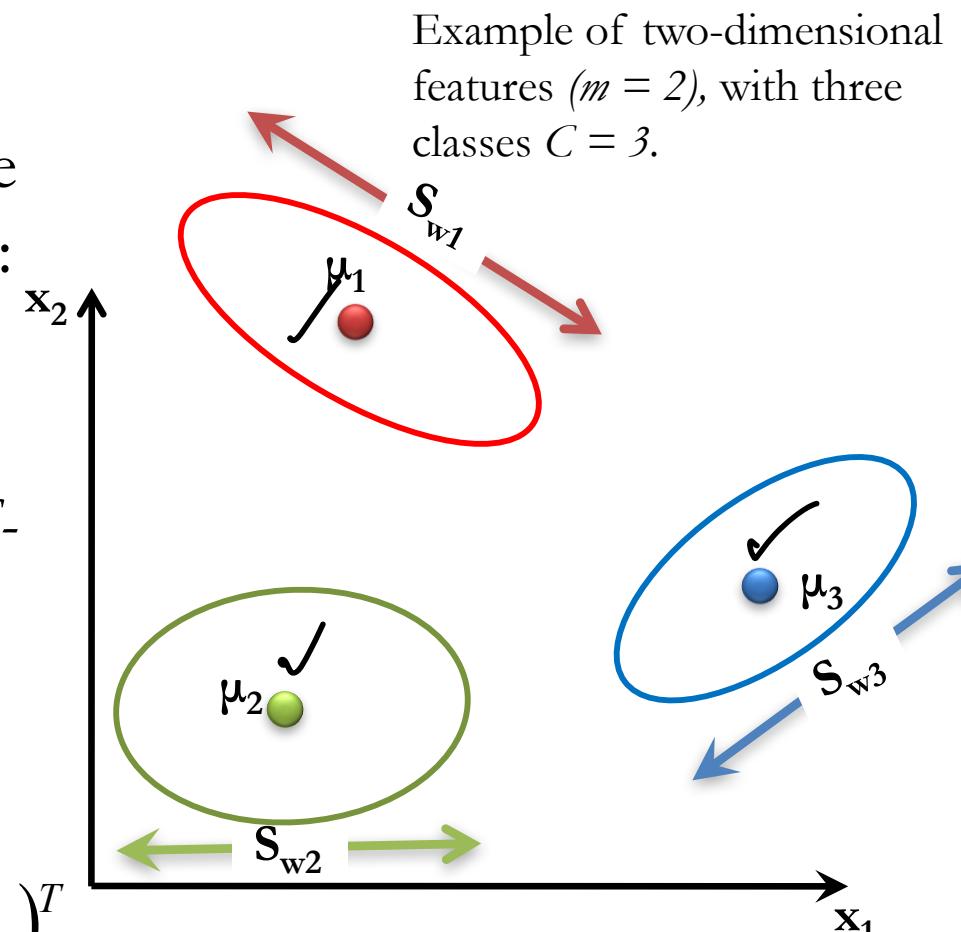
$$\underline{\underline{S_w}} = \underline{\underline{S_1}} + \underline{\underline{S_2}}$$

- This can be generalized in the C -classes case as:

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

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

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



N_i : number of data samples in class ω_i .

LDA – C-Classes

- Recall the two classes case, the *between-class scatter* was computed as:

$$\cancel{S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T}$$

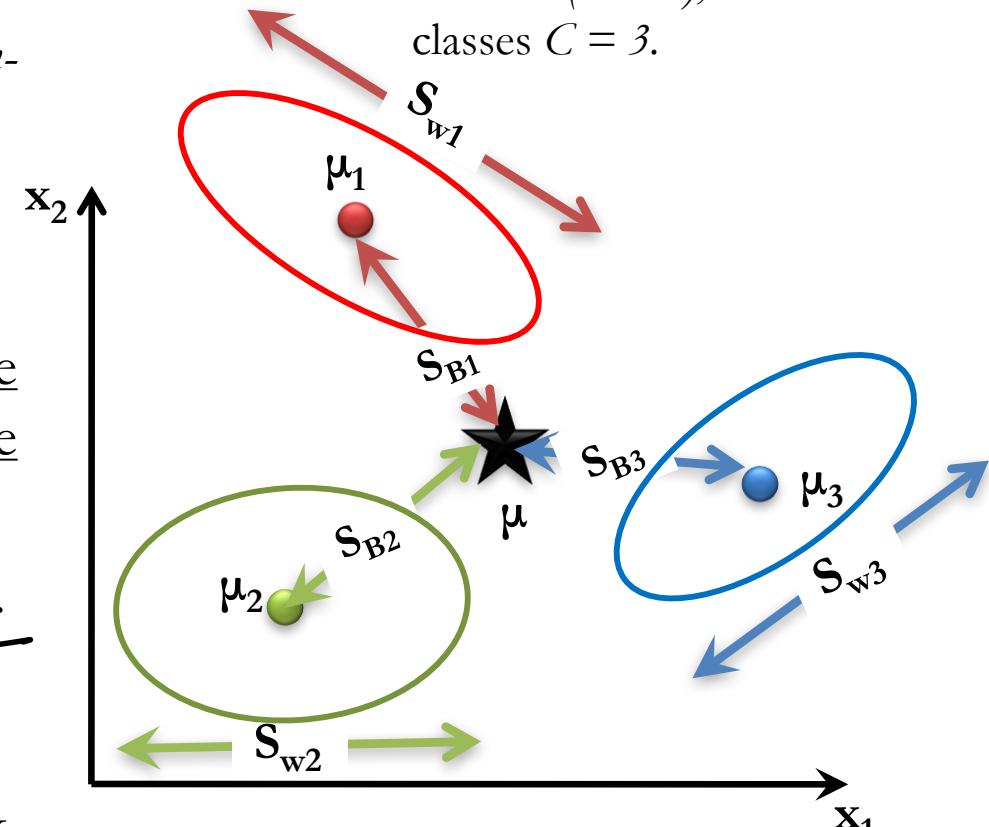
- For C -classes case, we will measure the between-class scatter with respect to the mean of all classes as follows:

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

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

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

Example of two-dimensional features ($m = 2$), with three classes $C = 3$.



N: number of all data .

N_i : number of data samples in class ω_i .

LDA – C-Classes

- Recall in two-classes case, we have expressed the scatter matrices of the projected samples in terms of those of the original samples as:

$$\Rightarrow \tilde{S}_W = W^T S_W W$$

$$\Rightarrow \tilde{S}_B = W^T S_B W$$

This still hold in C -classes case.

- Recall that we are looking for a projection that maximizes the ratio of between-class to within-class scatter.
- Since the projection is no longer a scalar (it has $C-1$ dimensions), we then use the determinant of the scatter matrices to obtain a scalar objective function:

$$J(W) = \frac{\left| \tilde{S}_B \right|}{\left| \tilde{S}_W \right|} = \frac{\left| W^T S_B W \right|}{\left| W^T S_W W \right|}$$

Matrix \mathbf{W} is $m \times n$
 \tilde{S}_B is $d \times d$
 \tilde{S}_W is $d \times d$
 The result is a 1×1 scalar

- And we will seek the projection \mathbf{W}^* that maximizes this ratio.

LDA – C-Classes

- To find the maximum of $J(W)$, we differentiate with respect to \mathbf{W} and equate to zero.
- Recall in two-classes case, we solved the eigen value problem.

$$\underline{S_W^{-1} S_B w = \lambda w} \quad \text{where } \lambda = J(w) = \text{scalar}$$

- For C -classes case, we have $C-1$ projection vectors, hence the eigen value problem can be generalized to the C -classes case as:

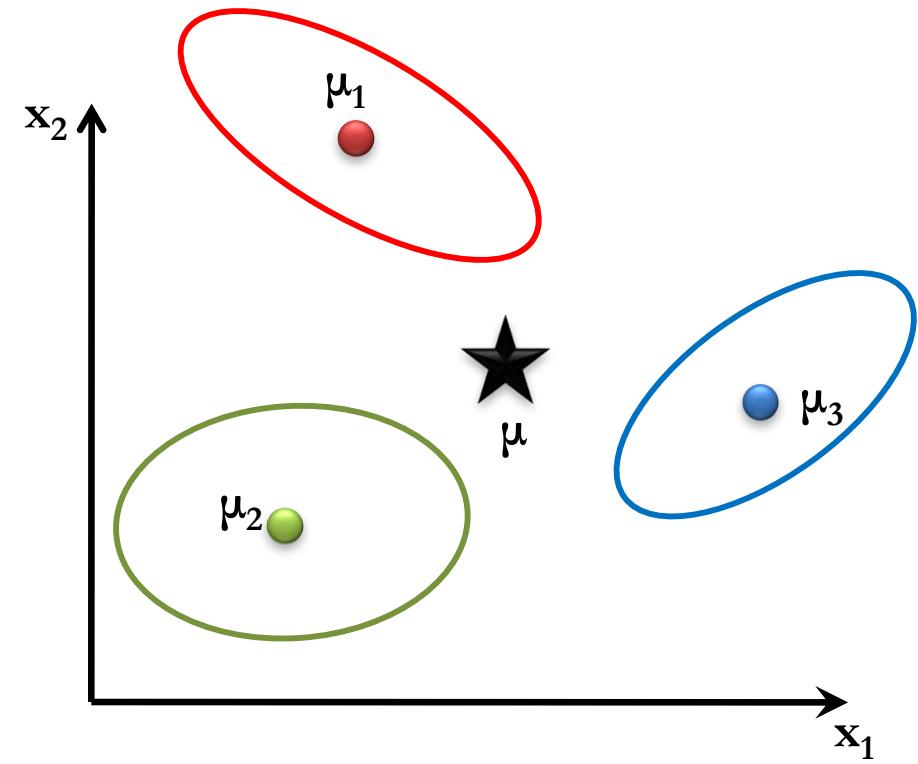
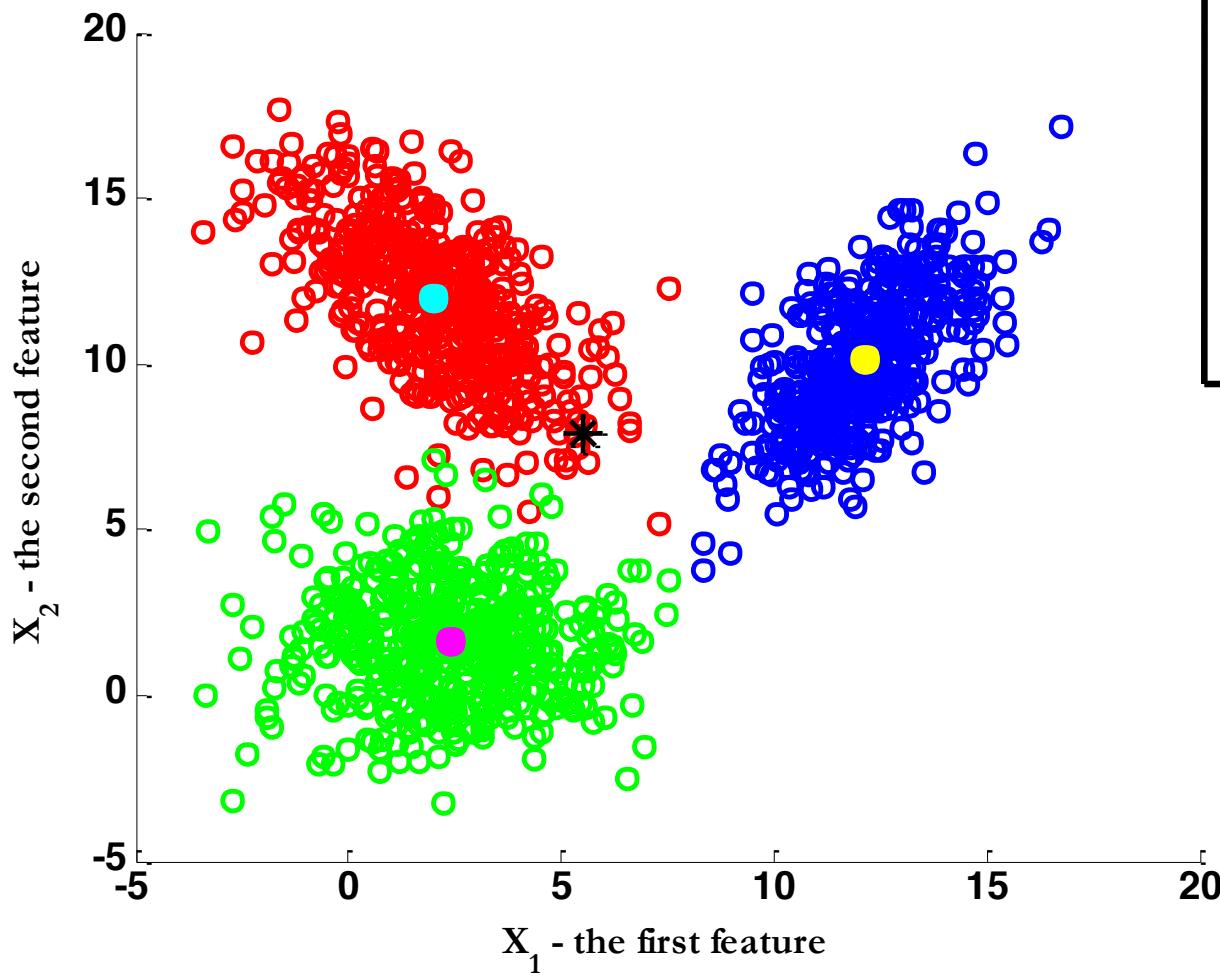
$$\underline{\underline{S_W^{-1} S_B w_i = \lambda_i w_i}} \quad \text{where } \underline{\underline{\lambda_i = J(w_i)}} = \text{scalar} \quad \text{and } i = 1, 2, \dots, C-1$$

- Thus, It can be shown that the optimal projection matrix \mathbf{W}^* is the one whose columns are the eigenvectors corresponding to the largest eigen values of the following generalized eigen value problem:

$$S_W^{-1} S_B W^* = \lambda W^*$$

where $\lambda = J(W^*) = \text{scalar}$ and $W^* = [w_1^* | w_2^* | \dots | w_{C-1}^*]$

It's Working ... 😊



Let's visualize the projection vectors W

```

%% 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

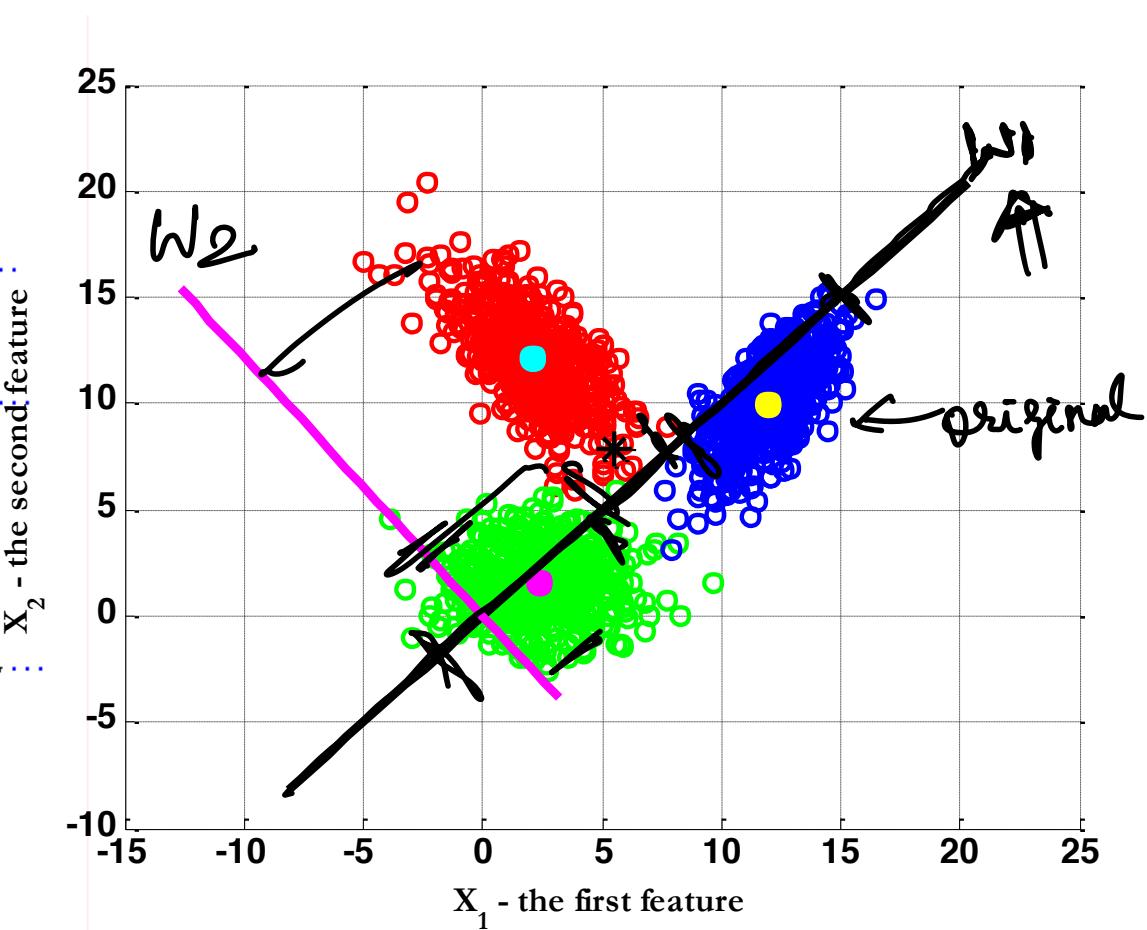
% the second class
scatter(X2(1,:),X2(2,:),'g','LineWidth',2,'Parent',axes1);
hold on

% 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(1);

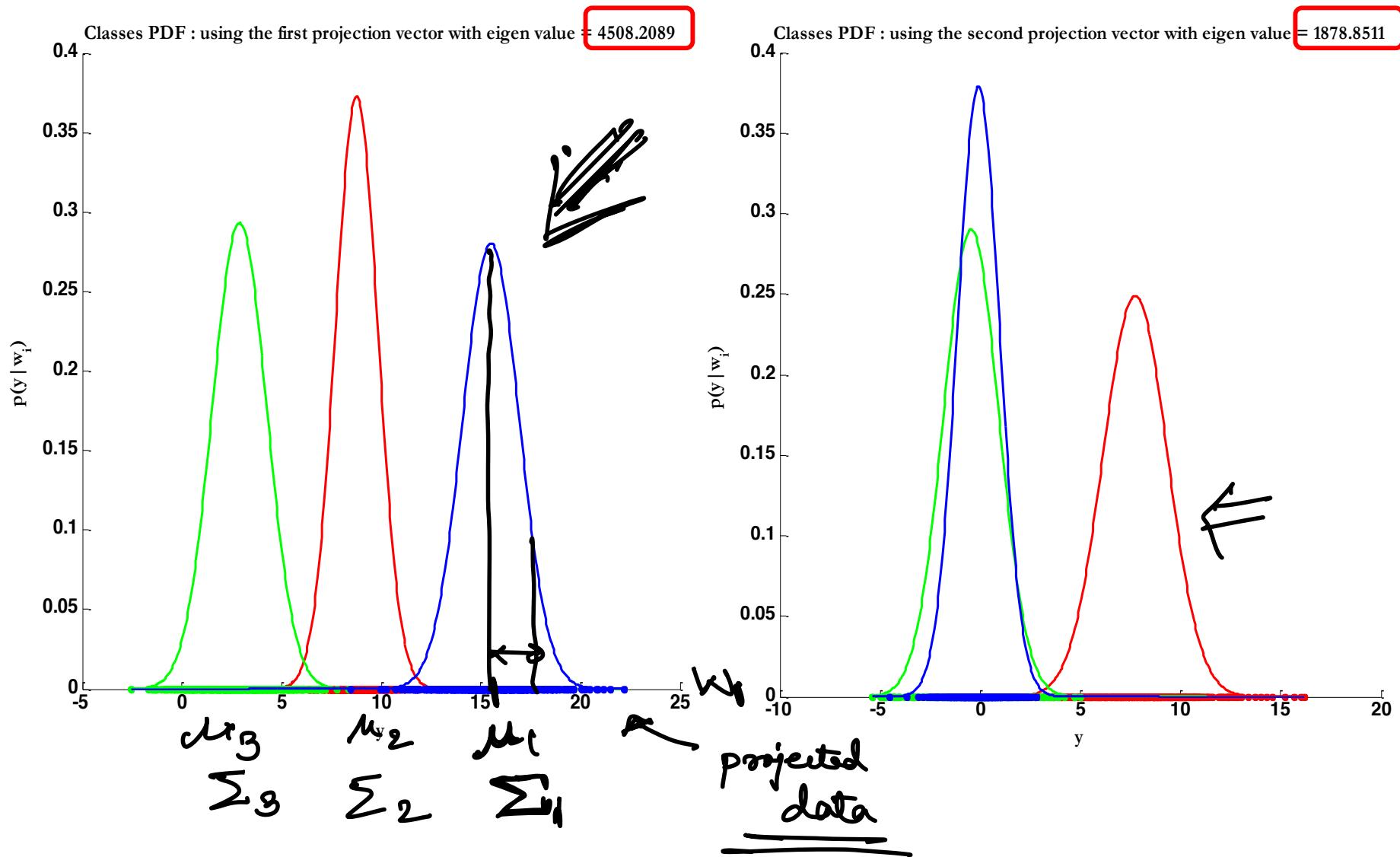
% the second vector
t = -5:20;
line_x2 = t .* w2(1);
line_y2 = t .* w2(2);

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

```

Which is Better?!!!

- Apparently, the projection vector that has the **highest eigen value** provides higher discrimination power between classes



$$N(x | \mu, \Sigma) = \frac{1}{(2\pi)^{D/2}} \frac{1}{|\Sigma|^{1/2}} \exp \left\{ -\frac{1}{2} (x - \mu)^T \Sigma^{-1} (x - \mu) \right\}$$

Dimension = $C-1$

any new test point

$$P(C_k | x) = P(x | C_k) \frac{P(C_k)}{\text{prior}}$$

x_{test} → $P_{C_1} | x_{\text{test}}$, μ_k, Σ_k

\vdots

$P_{C_m} | x_{\text{test}}$

$\left. \begin{array}{l} \\ \\ \end{array} \right\} \text{argmax} \frac{\# \text{ training points} \text{ in class } k}{\# \text{ training points}}$