## LDA : Linear Fisher Discriminator - Derivation

LDA is a special / specific case of MDA.

- only two classes
- and the approach to separate two classes is called Fisher Discriminator: also called as linear Fisher Discriminator.

## Fisher Discriminator :

## Given :

$n$    $d$-dimensional feature vectors

$$x_1 = (x_{11}, \dots x_{1d})$$
$$x_2 = (x_{21}, \dots x_{2d})$$
$$\vdots$$
$$x_n = (x_{n1}, \dots x_{nd})$$

out of these

$n_1$ no. of F.V $\in \omega_1$

$n_2$ no. of F.V $\in \omega_2$

$n_1 =$ no. of F.V from class $\omega_1$

$n_2 =$ no. of F.V from class $\omega_2$.

$\omega =$ direction of projection.

$\|\omega\| = 1$   (unit vector in the direction of projection line)

$\omega_1 = (x_1, x_2 \dots, x_{n_1}) = n_1$ no of samples

$\omega_2 = (x_1, x_2 \dots, x_{n_2}) = n_2$ no. of samples

Suppose If a take orthogonal projection on data $x_i$ and I get $y_i$

$$\boxed{y_i = w^t x_i}$$

$\Downarrow$

$$\boxed{y_i = w^t x_i}$$

$[w_1 \ w_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$w_1 x_1 + w_2 x_2 = 0$ and this passes through origin.
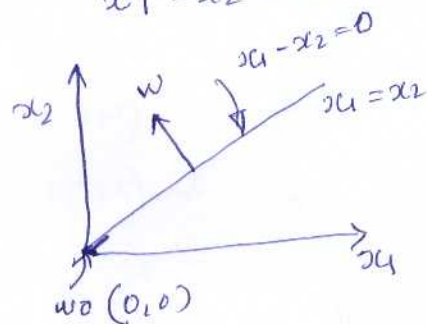
$x_1 = x_2$

$x_1 - x_2 = 0$

$w_1 x_1 - w_2 x_2 = 0$

$1 \cdot x_1 - 1 x_2 = 0$

$\underset{w^t}{(1 \quad -1)} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = 0$

$x_1 - x_2 = 0$



If I take projections on d-dimensional F.V $x_1, \ldots, x_n$

$y_1, y_2 \ldots , y_n$ [ these are projected vectors].

$$\underbrace{y_1, y_2, \ldots \ldots \ldots, y_n}_{\substack{n_1 \in \omega_1 \qquad n_2 \in \omega_2}}$$

let $m_1$ & $m_2$ is the mean of class $\omega_1$, $\forall x \in \omega_1$ and class

$\omega_2$ $\forall x \in \omega_2$

$$\boxed{m_1 = \frac{1}{n_1} \sum_{\forall x \in \omega_1} x}$$

$$\left[ \text{set } x = \underbrace{(x_1, x_2 \cdots x_{n_1})}_{n_1, \text{ no. of samples.}} \right] \text{ and}$$

$$\boxed{m_2 = \frac{1}{n_2} \sum_{\forall x \in \omega_2} x}$$ 

$\text{set } x = \{x_1, x_2 \cdots x_{n_2}\}$

$n_2$. no. of samples

$$\omega_1 = (x_1, x_2, \cdots x_{n_1})$$

$$\omega_2 = (x_1, x_{2_1}, \cdots x_{n_2})$$

when I compute the projection of the mean vector,

the projected mean vector $m_1$ can be written as

$\tilde{m}_1$.

$$\tilde{m}_1 = \frac{1}{n_1} \sum_{\forall y \in \omega_1} y \quad \leftarrow \text{ projected vector of } x.$$

$$y = w^t x.$$

$$\tilde{m}_1 = \frac{1}{n_1} \sum_{\forall x \in \omega_1} w^t x$$

$$= \frac{1}{n_1} w^t \sum_{\forall x \in \omega_1} x.$$

$$= w^t \frac{1}{n_1} \sum_{\forall x \in \omega_1} x = w^t m_1$$

$$\boxed{\tilde{m}_1 = w^t m_1}$$

$\tilde{m_1}$ ⇒ Projection of mean of vectors of Class set $\omega_1$ is nothing but $\omega^t m_1$.

In the same mannar, I can compute, the projection of mean of samples of class $\omega_2$ which is nothing but

$$\boxed{\tilde{m_2} = \omega^t m_2}$$

Now, what is the Distance between these two projected mean?

**Distance between projected means:**

$$|\tilde{m_1} - \tilde{m_2}| = |\omega^t m_1 - \omega^t m_2|$$

$$= |\omega^t (m_1 - m_2)|$$

vector $\omega$.

from, this I can ~~simply~~ increase the distance between $|\tilde{m_1} - \tilde{m_2}|$ by simply ~~scaling~~ scaling of '$\underline{\omega}$'.

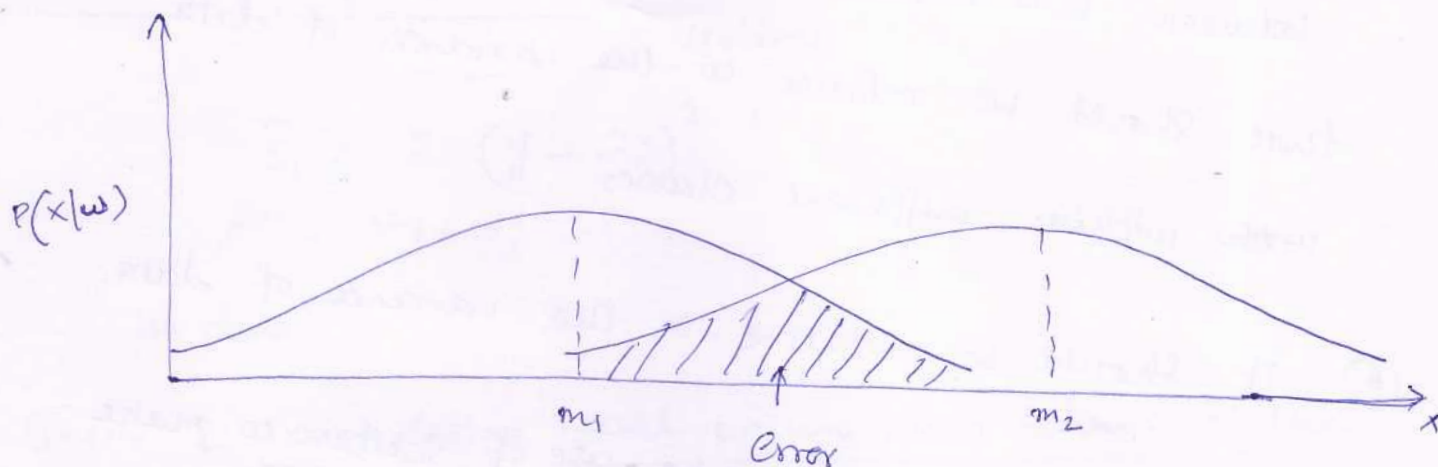Here I assume $\|\omega\| = 1$. If the $\|\omega\| > 1$ (more than 1) then the distance between $|\tilde{m_1} - \tilde{m_2}|$ will go on increasing. | * This ensures that I can increase the separability between two classes |

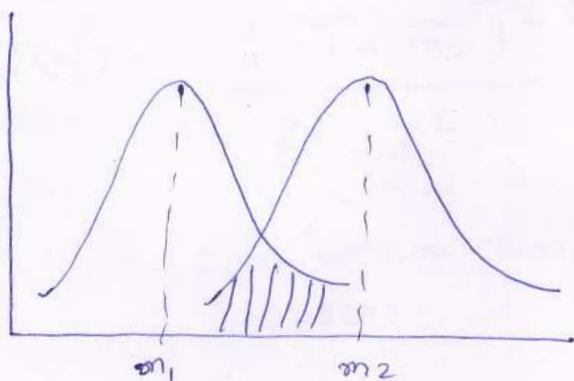But, the question is how much should I increase to separate the two classes?.

Illustration:

If the variance is large:



We should have the difference between $m_1$ and $m_2$ should be quite large so that the error of classification is reduced.

If the variance is small:



① In this case, we don't need that much separation which we need in the above case.

② The separation between $m_1$ and $m_2$ is much less than the above case case. (even then the error is minimized)

⊛ So how much should be the difference between two projected mean in the reduced space, that should be relative to the variance of data with within different classes.

⊛ It should be relative to the variance of data.

⊛ Accordingly, we can make use of <u>Scatter</u> to make some criterian function.

<u>Scatter of projected data:</u>

$$\tilde{s}_i^2 = \overbrace{\sum (y - \tilde{m}_i)^2}^{\text{variance}} \Rightarrow \text{we have not normalized } \frac{1}{n}.$$

$$\forall y \in \omega_i$$

↗ ith class

# Fisher Discriminant Method:

## Scatter of Projected data:

$$\tilde{S}_i = \overbrace{\sum (\tilde{y} - \tilde{m}_i)^2}^{\text{Variance.}}$$

$\forall \tilde{y} \in \omega_i$

↑ ith class

Then, I can define total within class scatter of the projected samples.

$$S = \tilde{S}_1^2 + \tilde{S}_2^2$$

We want the distance between two mean should be relative to this total scatter 's'. So we can define the criteria function as follows.

$$J(\omega) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$$

difference b/w projected mean ⟶ as large as possible w.r.t the total within class scatter.

⎵ total within class Scatter.

we want to maximize $J(\omega)$ which is the ratio of $|\tilde{m}_1 - \tilde{m}_2|^2$ upon $\tilde{S}_1^2 + \tilde{S}_2^2$. $J(\omega)$ is ratio of inter class scatter to intra class scatter.

$$J(\omega) = \frac{\text{Inter class}}{\text{Intra class}} = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$$

Fisher Discriminant maximize the

$$J(w) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

The value of 'w' which maximizes $J(w)$ is the projection direction.

### $J(\cdot)$ in terms of 'w'

$$S_i \uparrow \qquad\qquad S_w \rightarrow \text{total within class scatter.}$$

Scatter within ith class

$$S_i = \sum (x - m_i)(x - m_i)^t$$
$$\forall x \in w_i$$

Scatter for individual class
(ith class)

we can define, total within class scatter as

$$S_w = \sum_{\forall i} S_i \qquad \text{here, we consider only two classes}$$

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

Scatter of projected samples:

$$\tilde{S}_i = \sum_{\forall y \in w_i} (y - \tilde{m}_i)^2$$

$$\tilde{S}_i^2 = \sum_{\forall x \in w_i} (w^t x_i - w^t m_i)^2 \qquad y_i = w^t x_i$$

by rearranging this,

$$\tilde{S}_i^2 = \sum_{\forall x \in w_i} w^t (x_i - m_i)(x_i - m_i)^t w$$

$$\tilde{S}_i^2 = \sum_{\forall x \in w_i} (w^t x - w^t m_i)^2$$

$$= \sum_{\forall x \in w_i} (w^t x - w^t m_i) \cdot (w^t x - w^t m_i)$$

$$= \sum_{\forall x \in w_i} w^t (x - m_i) \cdot w^t (x - m_i)$$

$$= \sum w^t (x - m_i) \cdot (x - m_i)^t w.$$

$$= w^t \left[ \sum_{\forall x \in w_i} (x - m_i)(x - m_i)^t. \right] \cdot w$$

$$\boxed{\tilde{S}_i^2 = w^t S_i w.}$$

then the sum of $\tilde{s}_1^2 + \tilde{s}_2^2$

$$\tilde{s}_1^2 + \tilde{s}_2^2 = w^t S_1 w + w^t S_2 w$$

$$= w^t (S_1 + S_2) w$$

$$= w^t S_w w .$$

total within class scatter [refer page ④]

In the same mannar,

### Separation of the projected means.

$$(\tilde{m}_1 - \tilde{m}_2)^2 = (w^t m_1 - w^t m_2)^2$$

$$= (w^t m_1 - w^t m_2) \cdot (w^t m_1 - w^t m_2)$$

$$= w^t (m_1 - m_2) \cdot w^t (m_1 - m_2)$$

$$= w^t (m_1 - m_2) \cdot (m_1 - m_2)^t w$$

row vector       coloum vector.

$$= w^t \underbrace{(m_1 - m_2)(m_1 - m_2)^t}_{} w \quad \text{it rank is}$$

$$\quad\quad\quad n \quad\quad \text{atmost } 1.$$

outer product of '2' vectors

$$= w^t (\underbrace{S_B}_{} w) \Rightarrow \text{is a vector in the direction}$$

$$\text{of } (m_1 - m_2)$$

between class scatter

---

**Observation:** $w^t \underbrace{S_B w}_{} \Rightarrow$ is a vector in the direction of

$$(m_1 - m_2)$$

---

In particular, for any $w$, $S_B \cdot S_B w$ is in the direction of $m_1 - m_2$, and $S_B$ is quite singular.

$$J(\omega) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

in terms of '$\omega$'

$$J(\omega) = \frac{\omega^t S_B \omega}{\omega^t S_W \omega} \quad \Rightarrow \quad \text{generalized Rayleigh coefficient.}$$

↓ between class scatter

↑ total within class scatter

[ maximize this ratio by varying vector '$\omega$' and '$\omega$' gives

the projection direction.

Take a derivative w.r.t '$\omega$' and equate it to zero.

$$J(\omega) = \frac{\omega^t S_B \omega}{\omega^t S_W \omega} = \frac{u}{v} = \frac{v \cdot u' - u \cdot v'}{v^2}$$

$$\frac{d}{d\omega} J(\omega) = \frac{\omega^t S_W \omega \cdot \frac{d}{d\omega}(\omega^t S_B \omega) - \omega^t S_B \omega \frac{d}{d\omega}(\omega^t S_W \omega)}{(\omega^t S_W \omega)^2}$$

$$0 = \frac{\omega^t S_W \omega [2 S_B \omega] - \omega^t S_B \omega [2 S_W \omega]}{(\omega^t S_W \omega)^2}$$

$$\omega^t S_W \omega [2 S_B \omega] = \omega^t S_B \omega [2 S_W \omega]$$

Divide $\omega^t S_W \omega$

$$2 S_B \omega = \boxed{\frac{\omega^t S_B \omega}{\omega^t S_W \omega}} \, 2 S_W \omega.$$

$$\cancel{w'} S_B w = \lambda \cancel{w'} S_w w$$

$$\boxed{S_B w = \lambda S_w w} \implies \text{generalized eigenvalue}$$

$$\text{Problem}$$

If $S_w$ is non-singular then

$S_w^{-1}$ exists.

Singular = If the determinant
matrix     of matrix

$$|A| = 0$$

then singular

$$\boxed{S_w^{-1} S_B w = \lambda w} \implies \text{eigenvalue}$$

$$\text{Problem}$$

$(m_1 - m_2)$ ← eigen vector of $S_w^{-1} S_B$

and $\lambda$ is the corresponding

eigen value.

$$A^{-1} = \frac{1}{|A|} [ \ ]$$

Inverse does
not exist.
for singular
matrix

It is unnecessary to solve for the eigen values & vectors
of $S_w^{-1} S_B$ due to the fact that $S_B w$ is always in the direction of $m_1 - m_2$.

$$(\text{in } \lambda w)$$

The scaling factor of 'w' is not important $[\text{in } \lambda w]$

$$\boxed{S_B \cancel{w} = (m_1 - m_2) \cancel{w}}$$ but the direction of 'w' is so

important.

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

$$S_w^{-1} K (m_1 - m_2) = \lambda w$$

$$S_w^{-1} \left(\frac{K}{\lambda}\right) (m_1 - m_2) = w$$

$$S_w^{-1} (K_1) (m_1 - m_2) = w$$

$$\boxed{w = S_w^{-1} (m_1 - m_2)}$$