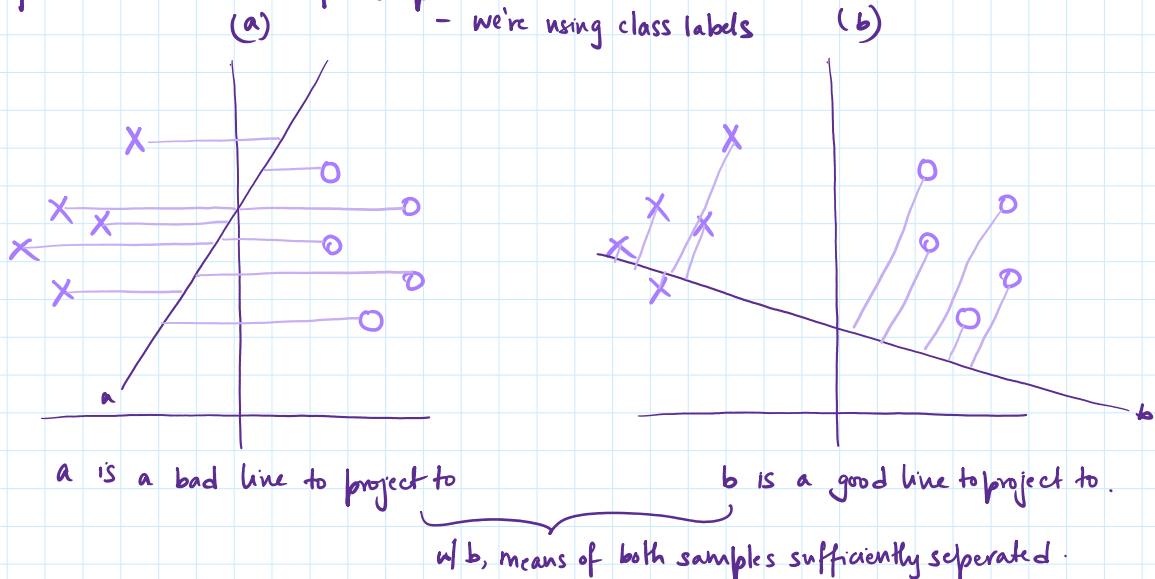


LDA for the same set of samples.



Deriving LDA.

Assumptions:

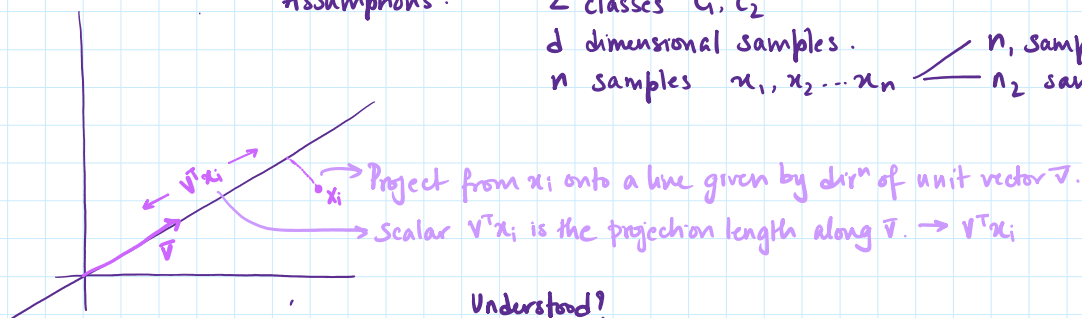
2 classes  $c_1, c_2$

d dimensional samples.

n samples  $x_1, x_2, \dots, x_n$

$n_1$  samples from  $c_1$

$n_2$  samples from  $c_2$



Understood?

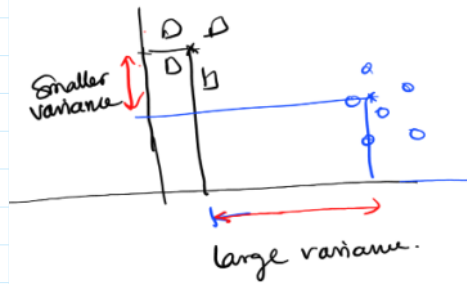
Let  $\bar{\mu}_1 + \bar{\mu}_2$  be the mean of projections of classes 1, 2.

$(\bar{\mu}_1 - \bar{\mu}_2)$  = measure of separation b/w classes.

$$\bar{\mu}_1 = \frac{1}{n_1} \sum_{x_i \in c_1} \vec{v}^T x_i = \vec{v}^T \left( \frac{1}{n_1} \sum_{x_i \in c_1} x_i \right) = \vec{v}^T \mu_1$$

$$\bar{\mu}_2 = \frac{1}{n_2} \sum_{x_i \in c_2} \vec{v}^T x_i = \vec{v}^T \left( \frac{1}{n_2} \sum_{x_i \in c_2} x_i \right) = \vec{v}^T \mu_2$$

$$\mu_2 = \frac{1}{n_2} \sum_{x_i \in C_2} v^T x_i = v^T \left( \frac{1}{n_2} \sum_{x_i \in C_2} x_i \right) = v^T \mu_2$$



But this doesn't look at the variance of the classes.

All we have to do is normalize  $|\mu_2 - \mu_1|$  — but how?

Keep that in mind.

Variance is our problem  
 $\Rightarrow$  normalize by a factor proportional to variance.

Scatter:

Say you're  $n$  samples:  $z_1, z_2, \dots, z_n$



Small scatter



Large scatter.

$\hookrightarrow$  Sample mean  $\mu_z = \frac{1}{n} \sum_{i=1}^n z_i$

Scatter  $\equiv$  variance BUT it's on a different scale.

$$\text{Scatter} = \sum (z_i - \mu_z)^2 \quad \text{or: (Sample variance)} \times n.$$

Connect the dots swappa

Normalizing  $|\mu_2 - \mu_1|$  by scatter = FISHER LDA. Linear discriminant analysis.

Let projected  $i$ th sample,  $y_i = v^T x_i$

Scatter for class 1  
 $\tilde{S}_1^2 = \sum (y_i - \bar{\mu}_1)^2$

Scatter for class 2  
 $\tilde{S}_2^2 = \sum (y_i - \bar{\mu}_2)^2$

Fisher LDA: "project on the line in the dir<sup>n</sup> of  $\bar{v}$  that MAXIMIZES  $J(\bar{v})$ "

$$J(\bar{v}) = \frac{(\bar{\mu}_1 - \bar{\mu}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2}$$

$\rightarrow$  If  $\tilde{S}_i^2$  is small, that means samples are crowded around  $\bar{\mu}_i$ .

Now express  $J$  explicitly as a function of  $v$ .  
 Then maximize!

① Define separate class matrices,  $S_1 + S_2$  for  $C_1 + C_2$  } scatter of original samples  $x_i$  BEFORE projection.

$$S_1 = \sum_{x_i \in C_1} (x_i - \mu_1)(x_i - \mu_1)^T$$

$$S_2 = \sum_{x_i \in C_2} (x_i - \mu_2)(x_i - \mu_2)^T$$

$\Rightarrow$  TOTAL scatter WITHIN classes  $S_W = S_1 + S_2$

② Remember  $\tilde{S}_1 + \tilde{S}_2$ ? Scatter for projected sample!

② Remember  $\tilde{S}_1 + \tilde{S}_2$ ? Scatter for projected sample!

$$\begin{aligned}\tilde{S}_1^2 &= \sum (y_i - \bar{u}_1)^2 \\ \tilde{S}_2^2 &= \sum (y_i - \bar{u}_2)^2\end{aligned} \quad \left. \vphantom{\sum} \right\} \text{look @ projection } y_i \text{ for a sample } x_i$$

$y_i = V^T x_i$   
 $\bar{u}_1 = V^T \mu_1$   
 $\bar{u}_2 = V^T \mu_2$

③ Plug in these values of  $y_i + \bar{u}_i$

$$\begin{aligned}\tilde{S}_1^2 &= \sum_{y_i \in c_1} (V^T x_i - V^T \mu_1)^2 \\ &= \sum [V^T (x_i - \mu_1)]^T [V^T (x_i - \mu_1)] \\ &= \sum \underbrace{[(x_i - \mu_1)^T V]^T}_{\text{transpose rule}} \underbrace{[V^T (x_i - \mu_1)]}_{\text{rearrange}} \\ &= \sum V^T (x_i - \mu_1) (x_i - \mu_1)^T V \\ \tilde{S}_1^2 &= V^T S_1 V \\ \Rightarrow \tilde{S}_2^2 &= V^T S_2 V\end{aligned}$$

④ Sum it all up now:

$$\begin{aligned}\tilde{S}_1^2 + \tilde{S}_2^2 &= V^T S_1 V + V^T S_2 V \\ &= V^T (S_1 + S_2) V \\ &= V^T S_W V !!!\end{aligned}$$

Quick review:

$$S_1 + S_2 = S_W$$

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

So we got scatter within,  $S_W$  (basically, scatter of data points within the class)



Between class scatter then becomes the difference b/w  $\bar{u}_i$  (means of classes)

Between class separation matrix  $S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T \rightarrow$  Before projection.



What about after projection?

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2$$

But we can write that better.

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = (V^T \mu_1 - V^T \mu_2)^2 \quad \text{cuz } \tilde{\mu}_i = V^T \mu_i$$

$$\begin{aligned}&= V^T (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T V \\ &= V^T S_B V\end{aligned}$$

Quick recap:

$$S_B = (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

$$(\tilde{\mu}_1 - \tilde{\mu}_2)^2 = V^T S_B V$$

Okay cool.

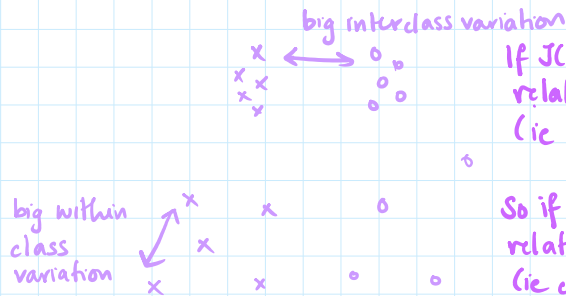
But why did we do all this?

$$\text{From before: } J(V) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2} = \frac{V^T S_B V}{V^T S_W V}$$

From before:  $J(V) = \frac{(\tilde{\mu}_1 - \tilde{\mu}_2)^2}{\tilde{S}_1^2 + \tilde{S}_2^2} = \frac{V^T S_B V}{V^T S_W V}$

## INTUITION TIME

$J(V) = \frac{V^T S_B V}{V^T S_W V} \rightarrow$  So  $J(V)$  ends up being a sort of quantitative measure for understanding variations b/w different classes as well as within.



If  $J(V) \uparrow$  it tells us that between class variation has increased relative to within class variation.  
(ie change in  $S_B \gg S_W$ )

So if  $J(V) \downarrow$  it tells us that within class variation has increased relative to b/w class variation  
(ie change in  $S_W \gg S_B$ )

Now for parameter estimation:

what we have:  $J(V)$

what we need to find: some value of  $v$  to optimise  $J(V)$

what really matters: NOT the magnitude of  $v$ . — remember

the general DIRECTION of  $v$ .  
cuz

$y_i = V^T x_i$ ?  
 $y$  is a  $f^n$  of  $x$ .  
 $\Rightarrow f(x_i) = V^T x_i$

Now before we go further, look @ this:

For a symmetric matrix  $W$  and a vector  $s$ , we have

$$\frac{\partial s^T W s}{\partial s} = 2 W s.$$

Notice that  $\Sigma_w$  is symmetric since its  $(i, j)^{th}$  element is

$$\sum_{n \in S_1} (x_{ni} - \mu_{1i})(x_{nj} - \mu_{1j}) + \sum_{n \in S_2} (x_{ni} - \mu_{2i})(x_{nj} - \mu_{2j}),$$

which is equivalent to its  $(j, i)^{th}$  element.

we want to rank  $f(x_i) = V^T x_i$  values.

That's all.

Using a vector proportional to  $v$  won't change ranking.

Done!

Now, 
$$\frac{d}{dv} J(V) = \frac{\left(\frac{d}{dv} V^T S_B V\right) V^T S_W V - \left(\frac{d}{dv} V^T S_W V\right) V^T S_B V}{(V^T S_W V)^2}$$

Now, we're trying to minimize  $J(V)$  so we'll set the derivative to 0.

Denom is a scalar so that goes.

$$0 = (2 S_B V) V^T S_W V - (2 S_W V) V^T S_B V$$

$$0 = S_B V V^T S_W V - S_W V V^T S_B V$$


$$0 = S_B V (V^T S_W V) - S_W V (V^T S_B V)$$

divide by  $V^T S_W V$  both sides.

$$0 = S_B V - \frac{(V^T S_B V)(S_W V)}{(V^T S_W V)} \rightarrow \lambda$$

$$\lambda = \frac{V^T S_B V}{V^T S_W V}$$

$$S_B V = \lambda S_W V$$


$$S_B v = \lambda S_W v$$

If  $S_W$  has a full rank (and  $\therefore$  is invertible)  
 $(S_W^{-1} S_B) v = \lambda v$

Just like we'd done w/ PCA, we obtain the LDA by taking