# 4.3 Linear Discriminant Analysis

## 4.3.1 Motivation.



feature mapping

Low-dimensional space

High-dimensional space

Observation: high dimensions make data classification easier, but data representation harder. It would be ideal if we can classify data in high dimensions then represent it in low dimensions.

Observation: PCA may fail to preserve
data separability.

Linear Discriminant Analysis (LDA): a method
to reduce dimensions while preserving
data separability.

Setting: sample pts with labels $(x^{(1)}, y_1) \cdots$
$(x^{(n)}, y_n)$ with $y_i = \pm 1$.
Set $\mathfrak{S}_1 \overset{def}{=} \{ x^{(i)} : y_i = 1 \}$

$$\mathcal{B}_2 \overset{\text{def}}{=} \{x^{(i)} : y_i = -1\}.$$

$$n_i \overset{\text{def}}{=} \# \text{ of sample pts in } \mathcal{B}_i . \quad i = 1, 2$$

Then the means of $\mathcal{B}_1$ and $\mathcal{B}_2$ are

$$\mu^{(1)} = \frac{1}{n_1} \sum_{x^{(i)} \in \mathcal{B}_1} x^{(i)}$$

$$\mu^{(2)} = \frac{1}{n_2} \sum_{x^{(i)} \in \mathcal{B}_2} x^{(i)}$$

The sample covariance matrices are

$$C^{(1)} = \frac{1}{n_1} \sum_{x^{(i)} \in \mathcal{B}_1} \left(x^{(i)} - \mu^{(1)}\right)\left(x^{(i)} - \mu^{(1)}\right)^T$$

$$C^{(2)} = \frac{1}{n_2} \sum_{x^{(i)} \in \mathcal{B}_2} \left(x^{(i)} - \mu^{(2)}\right)\left(x^{(2)} - \mu^{(2)}\right)^T$$

( If $\mathcal{B}_1$ is centered, then $\mu^{(1)} = 0$

and $\quad C^{(1)} = \frac{1}{n_1} X^{(1)} X^{(1)T}$ where

$$\mathbf{X}^{(1)} = \begin{pmatrix} | \\ x^{(i)} \\ | \end{pmatrix}_{x^{(i)} \in \mathcal{E}_1} \quad \text{is the sample matrix)}$$

Recall that $C^{(1)}$, $C^{(2)}$ are sym, pos semi-def.

First, we project $\mathcal{E}_1$, $\mathcal{E}_2$ to 1D.

Let $u$ be a unit vector, then the projected sample means are

$$\widetilde{\mu}^{(1)} \overset{\text{def}}{=} \text{Proj}_{\text{span}\{u\}} \mu^{(1)} = \left( \mu^{(1)T} u \right) u$$

$$\widetilde{\mu}^{(2)} \overset{\text{def}}{=} \text{Proj}_{\text{span}\{u\}} \mu^{(2)} = \left( \mu^{(2)T} u \right) u.$$
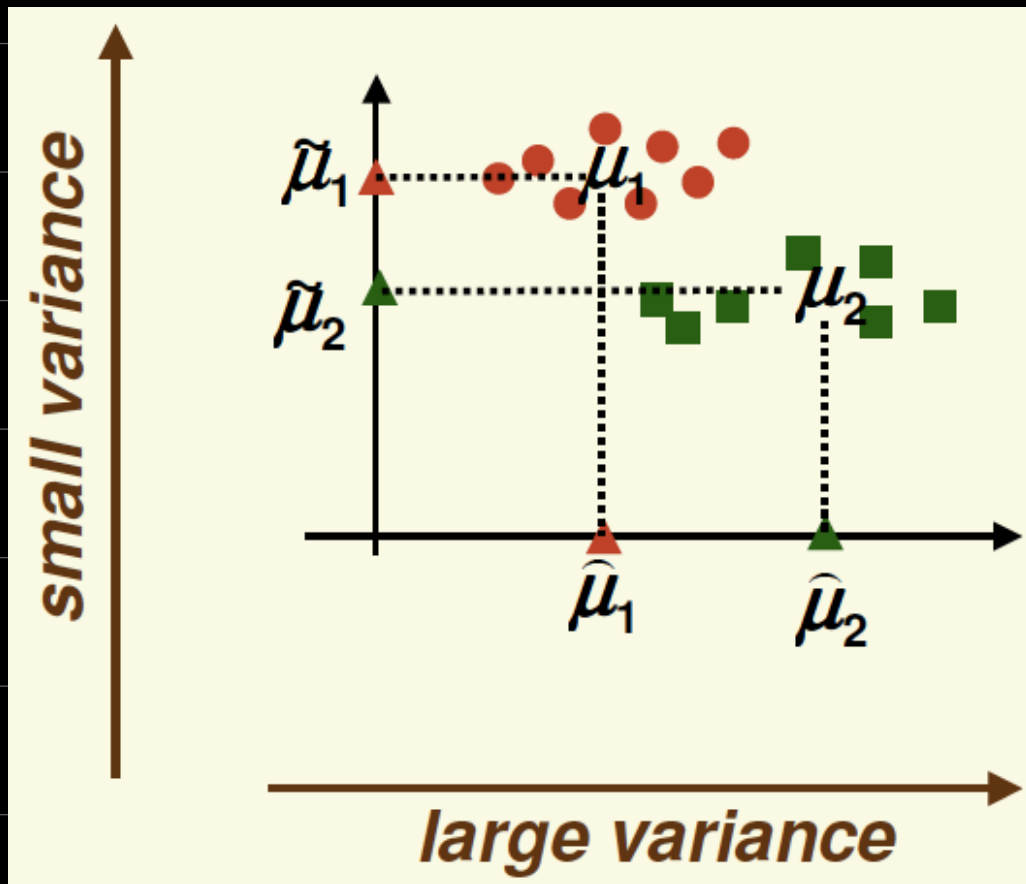
The projected variances are

$$\widetilde{C}^{(1)} = u^T C^{(1)} u$$

$$\widetilde{C}^{(2)} = u^T C^{(2)} u$$

(see Lecture 12)

## 4.3.2 Theory.

Q: What is a good metric for data separability?



Observation: $\|\tilde{\mu}^{(1)} - \tilde{\mu}^{(2)}\|$ itself is not a good metric, as data in each class may be scattered.

Fisher's idea: find a unit vector $u$ s.t.

$$\max_{\|u\|=1} \frac{\|\tilde{\mu}^{(1)} - x^{(2)}\|^2}{n_1 \tilde{C}_1 + n_2 \tilde{C}_2}$$

i.e. look for large $\|\tilde{u}^{(1)} - \tilde{u}^{(2)}\|$ with small scattering.

Notice: $\dfrac{\|\tilde{u}^{(1)} - \tilde{u}^{(2)}\|^2}{n_1 \tilde{C}^{(1)} + n_2 \tilde{C}^{(2)}}$ $\underset{\text{def of } \tilde{u}^{(i)} \text{ and } \tilde{C}^{(i)}}{=\!=\!=}$

$\dfrac{\|(u^{(1)T} u - u^{(2)T} u) u\|^2}{n_1 u^T C^{(1)} u + n_2 u^T C^{(2)} u}$ $\underset{\text{def of inner product}}{=}$

$\dfrac{u^T \overbrace{(u^{(1)} - u^{(2)})(u^{(1)} - u^{(2)})^T}^{S_B} u}{u^T \underbrace{(n_1 C_1 + n_2 C_2)}_{S_W} u} =$

$\dfrac{u^T S_B u}{u^T S_W u}$   $S_B, S_W$ are sym.