

Lecture 18

$$\text{LDA} \rightarrow \Pr\{Y=k \mid X=x\}$$

x

$$\hat{f}_k(x)$$

$$\hat{\pi}_k$$

$$\Pr\{X=x\}$$

$$\hat{\pi}_k = \frac{\#\{i: y_i = k\}}{n}$$

(an estimate
for $\Pr\{Y=k\}$)

$$\frac{1}{\sqrt{2\pi}\hat{\sigma}} e^{-\frac{(x - \hat{\mu}_k)^2}{2\hat{\sigma}^2}}$$

(an estimate
for
 $\Pr\{X=x \mid Y=k\}$)

At prediction time, given a new input

x :

$$\hat{y}(x) = \underset{k}{\operatorname{argmax}} \Pr\{Y=k \mid X=x\}$$

$$= \arg \max_k \frac{\hat{f}_k(x) \cdot \hat{\pi}_k}{\underbrace{p(x=x)}_{\text{the same for every class } k}}$$

$$= \arg \max_k \hat{f}_k(x) \hat{\pi}_k$$

$$\frac{1}{\sqrt{2\pi}\hat{\sigma}} e^{-\frac{(x-\hat{\mu}_k)^2}{2\hat{\sigma}^2}} = \arg \max_k \left\{ \log \hat{f}_k(x) + \log \hat{\pi}_k \right\}$$

$$= \arg \max_k \left\{ \log \frac{1}{\sqrt{2\pi}\hat{\sigma}} - \frac{(x-\hat{\mu}_k)^2}{2\hat{\sigma}^2} + \log \hat{\pi}_k \right\}$$

Because of the assumption that the variances were the same over all the classes, we can remove this quadratic term $-\frac{x^2}{2\hat{\sigma}^2}$

$$= \arg \max_k \left\{ \cancel{\log \frac{1}{\sqrt{2\pi}\hat{\sigma}}} + \cancel{\frac{x^2}{2\hat{\sigma}^2}} - \frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \frac{x\hat{\mu}_k}{\hat{\sigma}^2} + \log \hat{\pi}_k \right\}$$

do not depend on the class k

$$= \arg \max_k \left\{ -\frac{\hat{\mu}_k^2}{2\hat{\sigma}^2} + \frac{x\hat{\mu}_k}{\hat{\sigma}^2} + \log \hat{\pi}_k \right\}$$

linear in x

- Since the decision rule is linearly dependent on x , this classifier is called "linear" discriminant analysis.

Extension to the multi-dimensional setting:

$$\text{Data} = \left\{ (x_i, y_i) \right\}_{i=1}^n \quad \left\{ \begin{array}{l} y_i \in \{1, \dots, K\} \\ x_i \in \mathbb{R}^p \end{array} \right.$$

$$\begin{aligned} & \Pr\{y=K \mid X=x\} \\ &= \frac{\underbrace{f_K(x) \sim \mathcal{N}(\bar{\mu}_K, \Sigma)}_{\text{the same over all the classes}} \cdot \underbrace{\pi_K}_{\text{the same over all the classes}}}{\Pr\{X=x\}} \end{aligned}$$

How

Multi-variate Gaussian Distribution

$$(\mathcal{N}(\mu, \Sigma))$$

\downarrow \mathbb{R}^p \downarrow $p \times p$ matrix

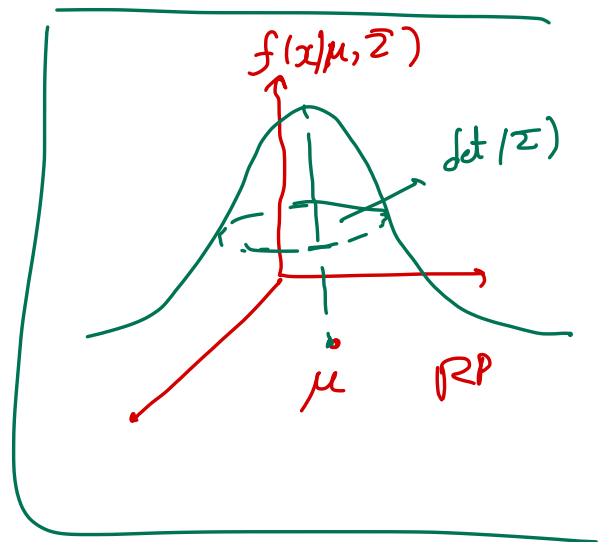
$$f(x | \overset{\text{mean}}{\mu}, \overset{\text{covariance matrix}}{\Sigma})$$

$$= \frac{1}{(2\pi)^{p/2} \cdot \det(\Sigma)^{1/2}} \exp \left(-\frac{1}{2} (x-\mu)^T \overset{\text{transpose}}{\Sigma^{-1}} (x-\mu) \right)$$

properties:

$$X \sim \mathcal{N}(\mu, \Sigma)$$

\downarrow
 $X = \begin{pmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{pmatrix} \in \mathbb{R}^p$



$$E[X] = \mu, \quad E[(X-\mu)(X-\mu)^T] = \Sigma$$

$\underbrace{\hspace{10em}}$
 $\forall r, s \in \{1, \dots, p\}$

$$E[(x_r - \mu_r)(x_s - \mu_s)] = \Sigma_{r,s}$$

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_r \\ \vdots \\ x_s \\ \vdots \\ x_p \end{pmatrix}$$

$$E[x_r] = \mu_r$$

$$E[x_s] = \mu_s$$

$$E[(x_r - \mu_r)(x_s - \mu_s)]$$

$$= \sum_{r,s}$$

$$\underbrace{\begin{pmatrix} & & & s \\ & & & \vdots \\ r & - & - & \sigma_{rs} \\ & & & \vdots \\ & & & s \end{pmatrix}}_{\text{matrix } \Sigma}^{p \times p}$$

$$\pi_k = \Pr\{Y=k\} \xrightarrow{\text{estimator}} \hat{\pi}_k = \frac{\overbrace{\#\{i: y_i=k\}}^{n_k}}{\text{total number of data points}}$$

$$f_k(x) = \mathcal{N}(\mu_k, \Sigma)$$

$$\hat{\mu}_k = \frac{\sum_{i: y_i=k} x_i}{n_k} \quad \text{number of data points with label } k$$

$$\hat{\Sigma} = \frac{1}{n-k} \begin{pmatrix} s & & \\ & \ddots & \\ & & \hat{\Sigma}_{rs} & \\ & & & \ddots \\ & & & & s \end{pmatrix}_{p \times p}$$

$$\hat{\Sigma}_{rs} = \frac{\sum_{k=1}^K \sum_{i: y_i = k} (X_{i,r} - \hat{\mu}_{k,r})(X_{i,s} - \hat{\mu}_{k,s})}{n-k}$$

we estimate $\hat{\mu}_k$ and $\hat{\Sigma}$ and $\hat{\tau}_k$ from data

\Rightarrow an estimate for the conditional probabilities

$$P\{y=k | X=x\} = \frac{P\{X=x | y=k\} P\{y=k\}}{P(X=x)}$$

estimated \Rightarrow

$$\frac{\hat{f}(x | \hat{\mu}_k, \hat{\Sigma}) \hat{\tau}_k}{P(X=x)}$$

ep.

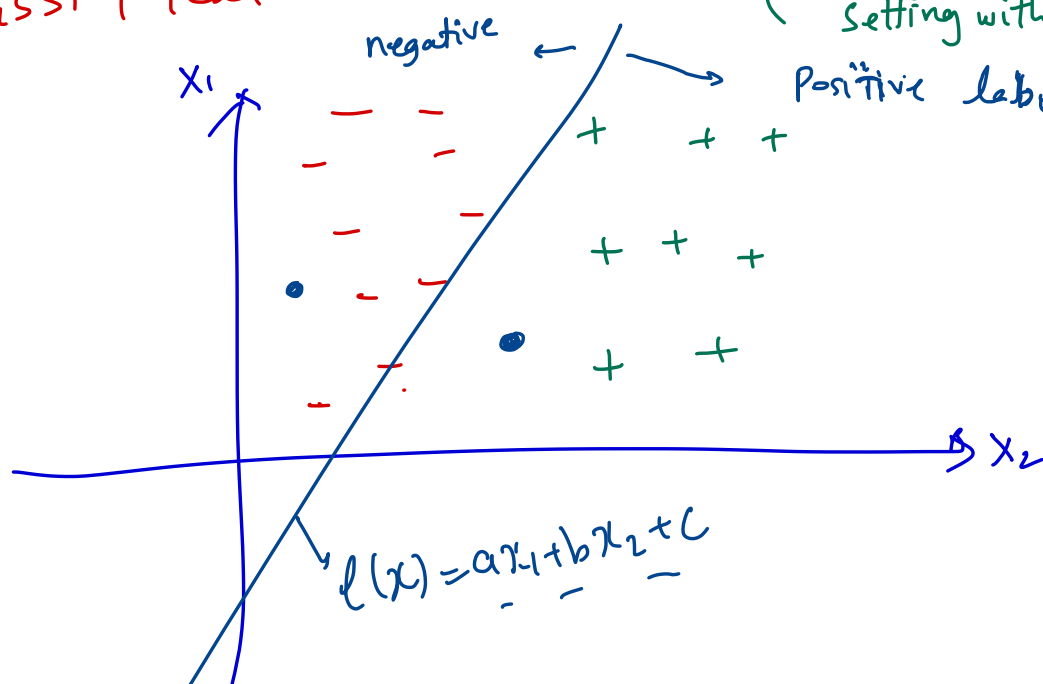
- Prediction of the label at a new data point x :

$$\hat{y}(x) = \underset{k}{\operatorname{argmax}} \left\{ \hat{f}(x | \hat{\mu}_k, \hat{\Sigma}) \hat{\pi}_k \right\}$$

$$= \underset{k}{\operatorname{argmax}} \left\{ \underbrace{x^T \hat{\Sigma}^{-1} \hat{\mu}_k - \frac{1}{2} \hat{\mu}_k^T \hat{\Sigma}^{-1} \hat{\mu}_k}_{\ell_k(x)} + \log \hat{\pi}_k \right\}$$

$\ell_k(x) \rightarrow \text{linear in } x$

This will lead to what we call linear classification boundaries. (let's consider a binary classification setting with 2-d data)



LDA: $\underset{k}{\operatorname{argmax}} \{ \ell_+(x), \ell_-(x) \}$

$$\Downarrow$$

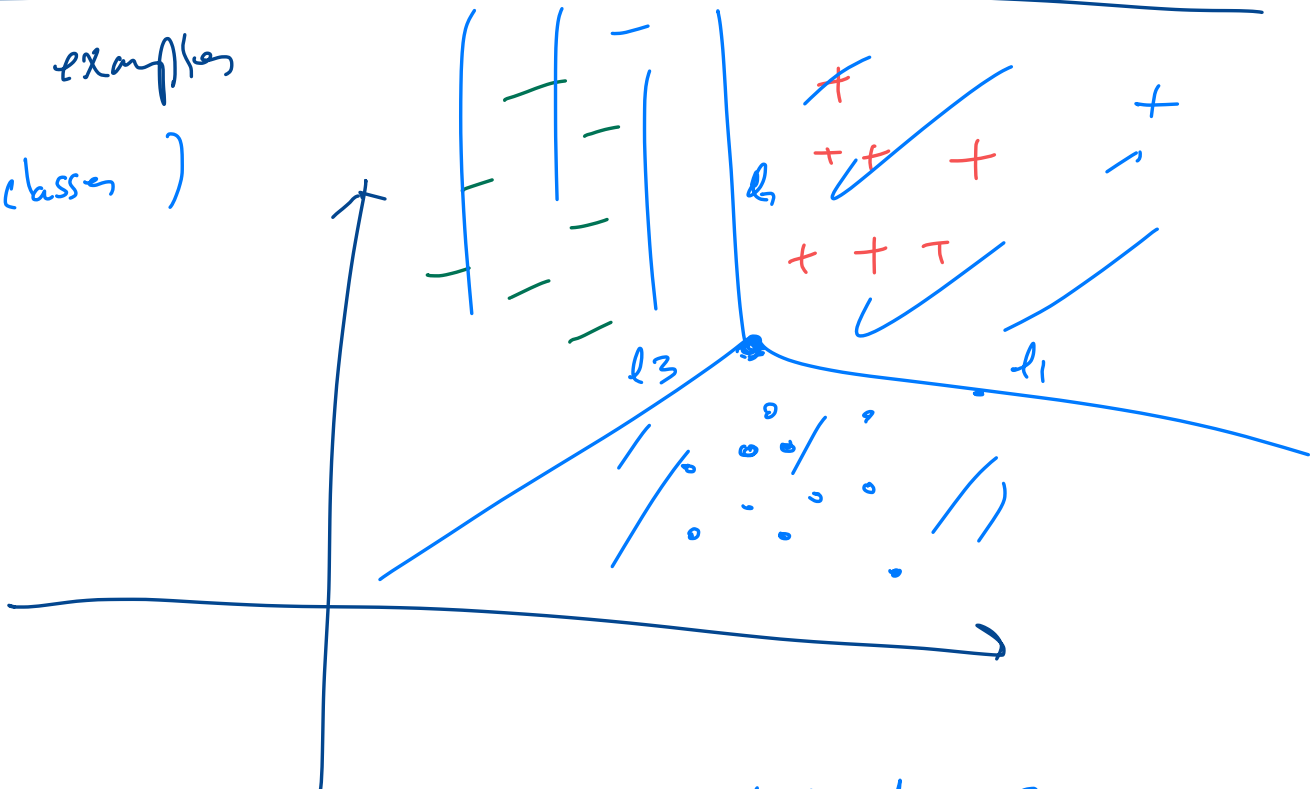
$$\left\{ \begin{array}{l} \text{if } l_+(x) > l_-(x) \Rightarrow \hat{y}(x) = + \\ \text{if } l_+(x) \leq l_-(x) \Rightarrow \hat{y}(x) = - \end{array} \right.$$

$$\Downarrow$$

$$\underbrace{l_+(x) - l_-(x)}_{\triangleq l(x)} \begin{cases} > 0 & \hat{y}(x) = + \\ < 0 & \hat{y}(x) = - \end{cases}$$

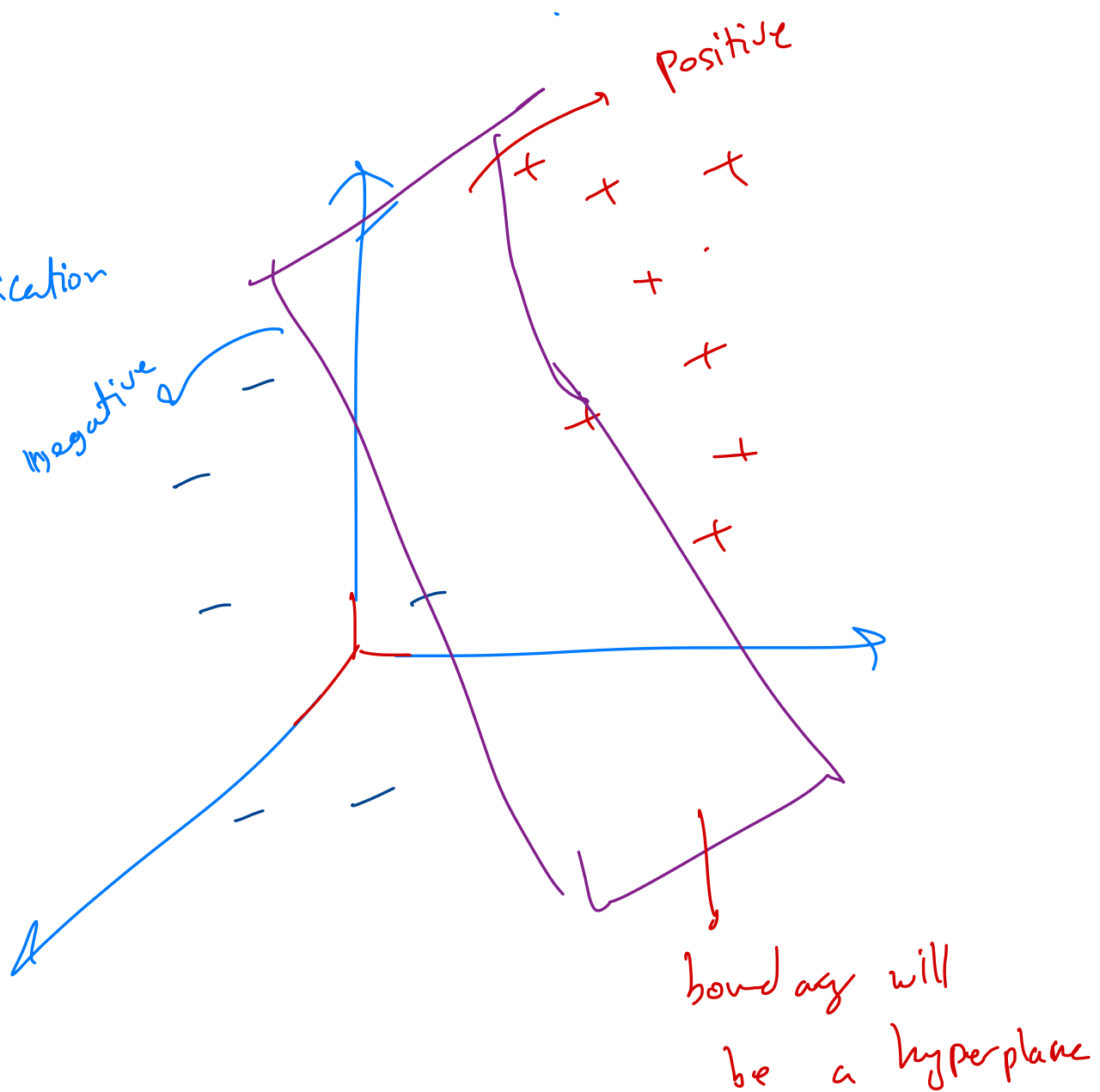
$$l(x) = ax_1 + bx_2 + c$$

other examples
(3 classes)

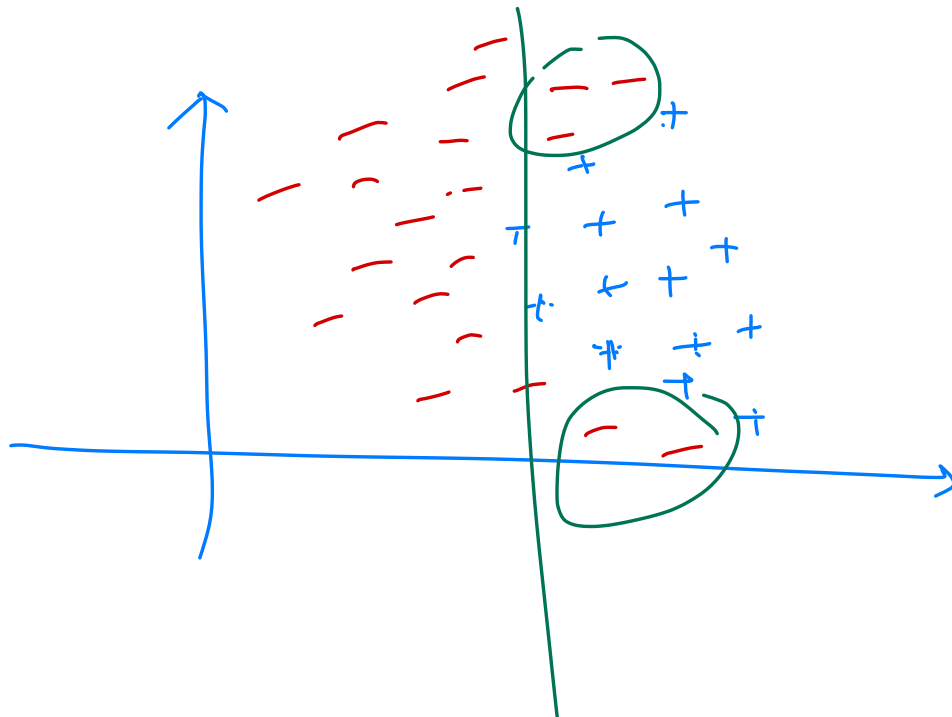


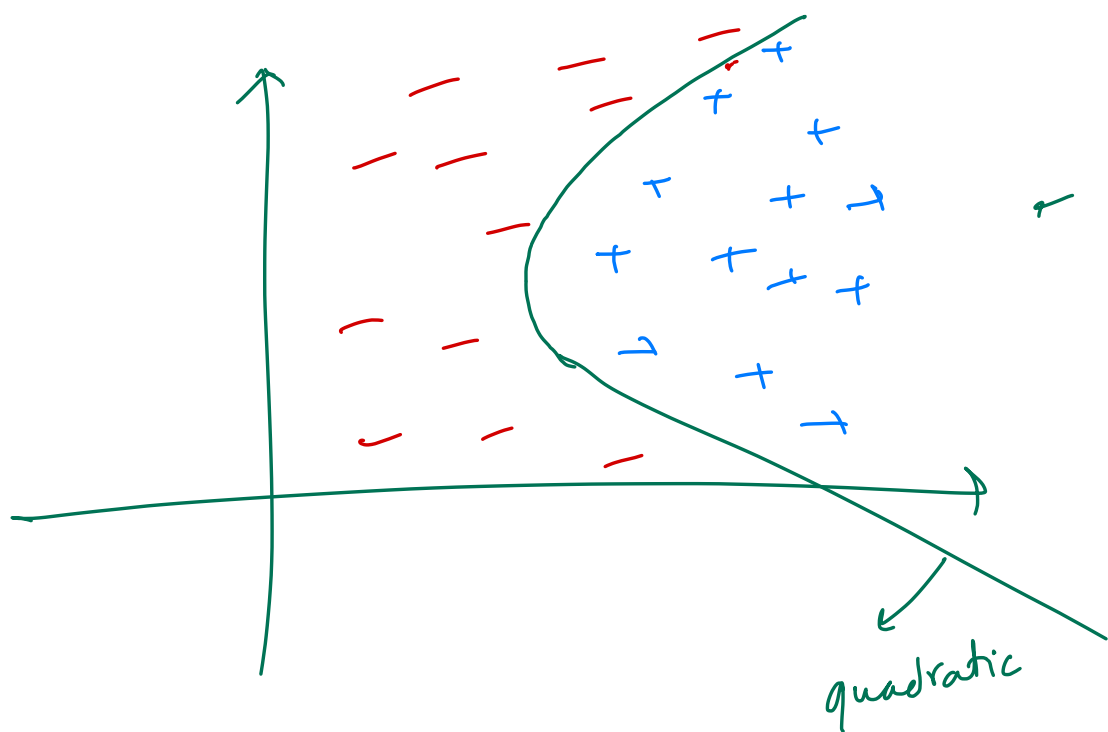
$$\hat{y}(x) = \operatorname{argmax} \{ l_+(x), l_-(x), l_0(x) \}$$

A
3-dim
binary
classification



In what cases can LDA be ineffective?





Quadratic Discriminant Analysis:

- Similar to LDA except that the variances can change across the classes.
 - for each class k : $f_k(x) = \mathcal{N}(\hat{\mu}_k, \hat{\Sigma}_k)$
 - learned using data from class k .
 - quadratic term
- $$\hat{y}_{QDA}(x) = \underset{k}{\operatorname{argmax}} \left\{ \underbrace{-\frac{1}{2}(x - \hat{\mu}_k)^T \hat{\Sigma}_k^{-1}(x - \hat{\mu}_k)}_{\text{quadratic term}} - \frac{1}{2} \log \det(\hat{\Sigma}_k) + \log \hat{\pi}_k \right\}$$

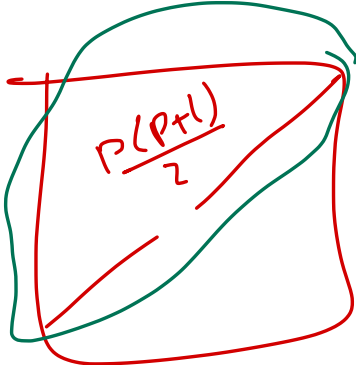
- Why QDA is a more complex class than LDA?

LDA

→ learn the mean vector per class + the covariance matrix.

$$K \cdot p + \frac{p(p+1)}{2}$$

means → covariance matrix → symmetric matrix

$\Sigma =$ 

QDA

- 1 mean vector per class
- 1 covariance matrix per class

$$K \cdot p + K \frac{p(p+1)}{2}$$

$$Kp + \frac{p(p+1)}{2}$$

< <
↓

when p is large

$$Kp + \frac{Kp(p+1)}{2}$$

LDA

→ less parameters

→ less flexibility / complexity

might underfit

QDA

might overfit

LDA is a better bet than QDA if there are relatively few training data points. In contrast, QDA is a better choice if the number of training data points is large or if the assumption of a common covariance matrix for all the k classes is off.

(in this case LDA will underfit)

Module 4 : Unsupervised learning:

Supervised learning: $\{ \underset{=}{x_i}, \underset{=}{y_i} \}_{i=1}^n \longrightarrow y_i \approx \bar{f}(x_i)$

unsupervised learning $\{x_i\}_{i=1}^n$