

- Labs are now online till the end of 2021
- we will use the same zoom link as for the lecture

LDA/QDA

model distribution of $X \in \mathbb{R}^p$ separately in each class ($k=1, 2, \dots, K$) and Bayes' theorem

$$P_r(Y=k | X=x) = \frac{\pi_k f_k(x)}{\sum_{l=1}^K \pi_l f_l(x)}$$

$$\pi_k = P_r(Y=k) \stackrel{\text{estimate}}{=} \frac{\# \{Y=k\}}{\# \{X\}}$$

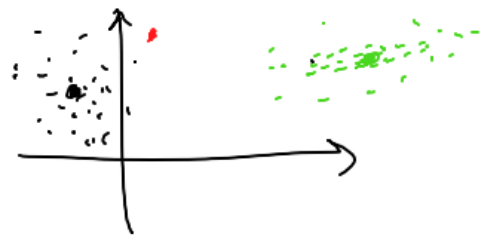
$f_k(x)$ - estimated density for

$$P(X \in A | Y=k) = \int_A f_k(x) dx$$

LDA/QDA: we assume that f_k are Gaussian densities

assume $p=1$ $f_k(x) = \frac{1}{\sqrt{2\pi} \sigma_k} \exp\left(-\frac{1}{2\sigma_k^2} (x - \mu_k)^2\right)$, $k=1, \dots, K$

In LDA we additionally assume that $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_k^2 = \sigma^2$
(but in QDA)



LDA, p=1

$$\text{model: } Pr(Y=k | X=x) = \frac{\pi_k \cdot \frac{1}{\sqrt{2\pi}\sigma} \cdot \exp\left(-\frac{1}{2\sigma^2} (x-\mu_k)^2\right)}{\sum_{k=1}^K \pi_k \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2\sigma^2} (x-\mu_k)^2\right)} \quad (*)$$

classification: $\text{argmax}_k Pr(Y=k | X=x)$

Since denominators in (*) are the same for all k, we may look at the numerator only, taking log:

$$\log \pi_k + \log \frac{1}{\sqrt{2\pi}\sigma} - \frac{1}{2\sigma^2} (x-\mu_k)^2 \quad \leftarrow \text{find k for which this is the largest}$$

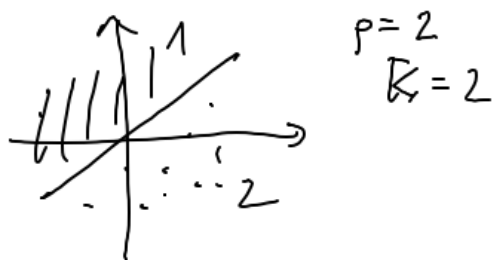
$$\log \pi_k + \frac{1}{\sigma^2} x \mu_k - \frac{1}{2\sigma^2} \mu_k^2 \quad \leftarrow \text{--- " ---}$$

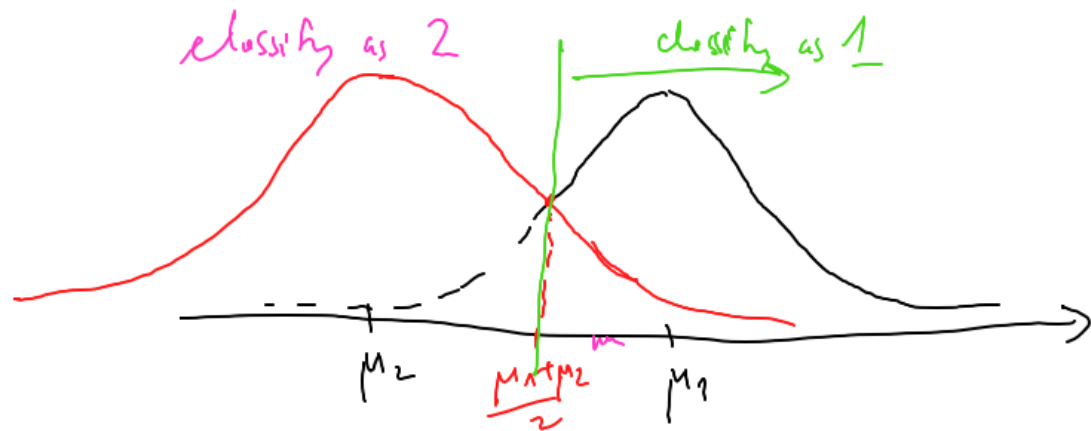
linear in x

assuming $k=2$, $\pi_1 = \pi_2$: $\frac{1}{\sigma^2} \mu_1 x - \frac{1}{2\sigma^2} \mu_1^2 > \frac{1}{\sigma^2} \mu_2 x - \frac{1}{2\sigma^2} \mu_2^2$

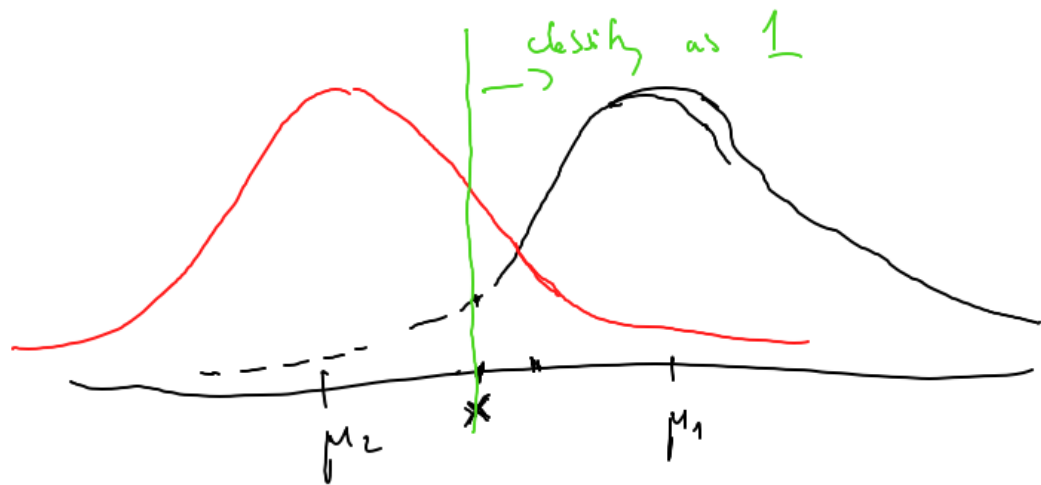
$\mu_1 > \mu_2$ $\frac{1}{\sigma^2} x (\mu_1 - \mu_2) > \frac{1}{2\sigma^2} (\mu_1^2 - \mu_2^2)$

$x > \frac{1}{2} (\mu_1 + \mu_2)$





$$\bar{\pi}_1 \geq \bar{\pi}_2$$



$$\pi_1 > \pi_2$$

How to obtain $\hat{\pi}_k, \hat{\mu}_k, \hat{\sigma}^2$:

$$\hat{\pi}_k = \frac{n_k}{n}$$

$$\hat{\mu}_k = \frac{1}{n_k} \sum_{i: y_i=k} x_i$$

$n_k = \#\{i: y_i=k\}$ - number of samples in class k

$$n = n_1 + n_2 + \dots + n_K$$

$$\hat{\sigma}^2 = \frac{1}{n-K} \sum_{k=1}^K \sum_{i: y_i=k} (x_i - \mu_k)^2$$

$p \gg 1$

$$f_k(x) = \frac{1}{(2\pi)^{p/2} |\det \Sigma|^{1/2}} \exp\left(-\frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k)\right)$$

$$\Sigma = \text{Cov} X = [\text{Cov}(X_i, X_j)]_{i,j}$$

$$\begin{aligned} \text{Cov}(X_i, X_j) &= \\ &= E((X_i - EX_i)(X_j - EX_j)) \end{aligned}$$

$$\leadsto f_k(x) = x^T \Sigma^{-1} \mu_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \log \pi_k$$

$\arg \max_k f_k(x) \leadsto$ prediction

$\hat{\pi}_k, \hat{\mu}_k$ - as before

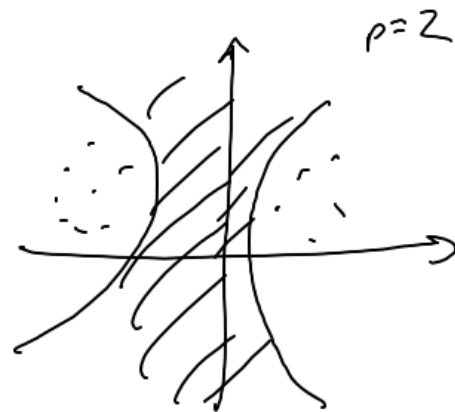
matrix $p \times p$

$$\hat{\Sigma} = \frac{1}{n-K} \sum_{k=1}^K \sum_{i: y_i=k} \begin{bmatrix} (x_i - \hat{\mu}_k) & (x_i - \hat{\mu}_k)^T \end{bmatrix}$$

QDA: as LDA, but Σ depends on the class: $\Sigma_k, k=1, \dots, K$

$$f_k(x) = -\frac{1}{2}(x - \mu_k)^T \cdot \Sigma_k^{-1}(x - \mu_k) - \frac{1}{2} \log |\det \Sigma_k| + \log \pi_k$$

$$= \underbrace{\sum a_j^{(k)} x_i \cdot x_j}_{\text{quadratic}} + \underbrace{b_j^{(k)} x_j}_{\text{linear}} + \underbrace{c^{(k)}}_{\text{const.}}$$



$$\left\{ \hat{\Sigma}_k = \frac{1}{n_k - 1} \sum_{i: y_i = k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T \right.$$

in QDA the number of parameters to estimate may be larger than

in LDA: $\Sigma \sim p \times p = p^2$

QDA: $\Sigma_k \sim p \times p$
 $\underbrace{K \cdot p^2}$

Naive Bayes

The assumption is that on each $\{Y=k\}$, X_1, X_2, \dots, X_p (the predictors) are independent
to classify an observation $x \in \mathbb{R}^p$:

$$k: \operatorname{argmax}_k \underbrace{P(Y=k)}_{\pi_k} \cdot \underbrace{P(X_1=x_1 | Y=k) \cdots P(X_p=x_p | Y=k)}$$

estimated e.g. by assuming some form of the density
and by estimating its parameters

Maximal margin classifier

{ binary classification

• Suppose that $(x_i, y_i) \in \mathbb{R}^p \times \{-1, 1\}$

- for ^(*) simplicity of interpretation let us assume that there exist $\beta_0 \in \mathbb{R}$ and $\beta_1 \in \mathbb{R}^p$ such that

$$\beta_0 + \beta_1 \cdot x_i > 0 \quad \text{if } y_i = 1$$

$$\beta_0 + \beta_1 \cdot x_i < 0 \quad \text{if } y_i = -1$$

scalar product

$\{x: \beta_0 + \beta_1 \cdot x = 0\}$ - hyperplane in \mathbb{R}^p

$(\beta_0 + \beta_{11} \cdot x_1 + \beta_{12} \cdot x_2 = 0 \quad \text{in } \mathbb{R}^2)$

(i.e., we assume that there exists a separating hyperplane)

so that

this makes sense without (*)

For the max. margin classifier we choose β_0, β_1

maximize M for which $\exists \beta_0, \beta_1$

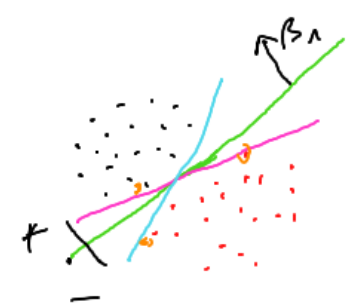
$$\min_i (\beta_0 + \beta_1 \cdot x_i) \cdot y_i \geq M$$

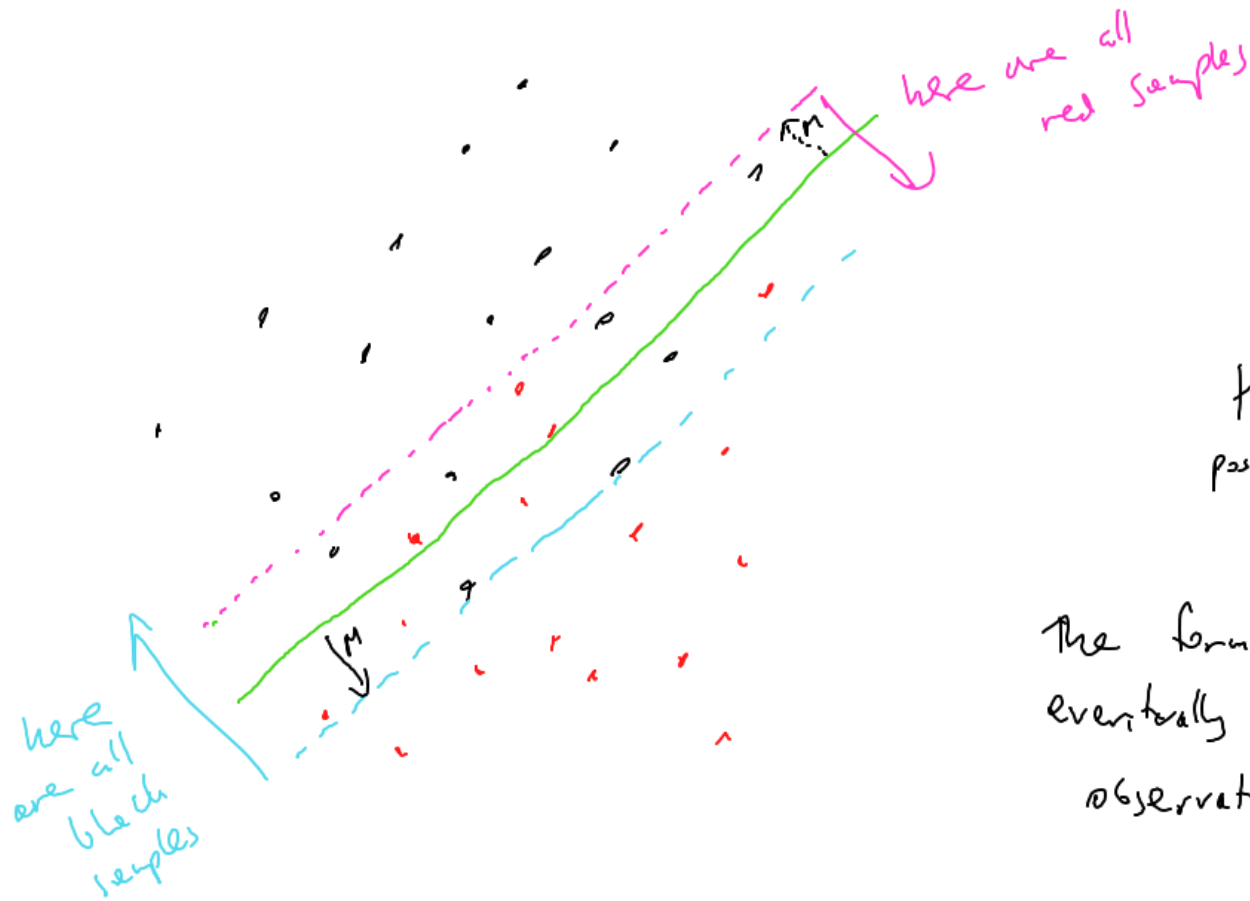
distance between x_i

and the hyperplane $\{x: \beta_0 + \beta_1 \cdot x = 0\}$

$$\|\beta_1\|_2 = 1$$

\rightarrow classification: $\text{sgn}(\beta_0 + \beta_1 \cdot x)$





here the maximum M
possible will be negative.

The formula for the classifier
eventually depends on just few
observations

