

LDA model $f_k(x) = \Pr(X=x | Y=k)$

$$\pi_k = \Pr(Y=k)$$

$$\Pr(X \in A | Y=k) = \int_A f_k(y) dy$$

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

classifier: classify x to the class k for which $\Pr(Y=k | X=x)$ is the largest

LDA: ($p=1$)

Take $f_k(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu_k)^2}{2\sigma^2}\right)$

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

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

decision boundary is a hyperplane



LDA for $p \geq 1$

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

$$\Sigma = \text{Cov} X = \begin{matrix} [\text{Cov}(X_i, X_j)]_{i,j} \\ \parallel \\ E[(X_i - E X_i)(X_j - E X_j)] \end{matrix}$$

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

\leadsto k with the largest $\delta_k(x)$
is the predicted class

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

④

QDA (Quadratic Discriminant Analysis)

like LDA, but with σ_k^2 or Σ_k depending on the class $k \in \{1, \dots, K\}$

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

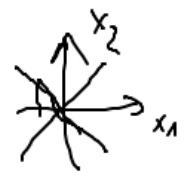
Classifier: take k for which $\delta_k(x)$ is the largest.

Ex: $p=1$

$$-\frac{1}{2} (x - \mu_k)^2 \sigma_k^{-2} - \frac{1}{2} \cdot 2 \log \sigma_k + \log \pi_k \geq -\frac{1}{2} (x - \mu_j)^2 \sigma_j^{-2} - \frac{1}{2} \cdot 2 \log \sigma_j + \log \pi_j$$

$$-\frac{1}{2} x^2 \sigma_k^{-2} + \dots \geq -\frac{1}{2} x^2 \sigma_j^{-2} + \dots$$

↑
if $\sigma_k \neq \sigma_j$, these do not cancel out



$$\left\{ \begin{array}{l} x_1^2 - x_2^2 = 0 \\ (x_1 - x_2)(x_1 + x_2) = 0 \end{array} \right.$$

$$\left\{ \begin{array}{l} x_1^2 + x_2^2 + 1 \geq 0 \\ -x_1^2 + x_2^2 - 1 \geq 0 \end{array} \right.$$



if $p=2, K=2$,
then the decision region
is an ellipse, a parabola, a hyperbola,
2 lines, line or \emptyset

Compared to LDA, QDA is a more flexible model with more parameters
(many more if p is large)

← LDA is yet another method that gives a linear decision boundary



Naive Bayes

the assumption is that on each $\{Y=c\}$, X_1, \dots, X_p (the predictors) are independent

$$k: \operatorname{argmax}_k P(Y=k) \cdot \underbrace{P(X_1=x_{01} | Y=k) \dots}_{\text{estimated e.g. by assuming some form of the density}} \cdot P(X_p=x_{0p} | Y=k)$$

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

Maximal margin classifier

Suppose that $(X_i, Y_i) \in \mathbb{R}^p \times \{-1, 1\}$ and suppose that the classes can be separated by a hyperplane, i.e., there exists $\beta_0 \in \mathbb{R}, \beta_1 \in \mathbb{R}^p$ such that $\beta_1 \neq (0, \dots, 0)$

$$\beta_0 + \beta_1 \cdot X_i > 0 \quad \text{if } Y_i = 1$$

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



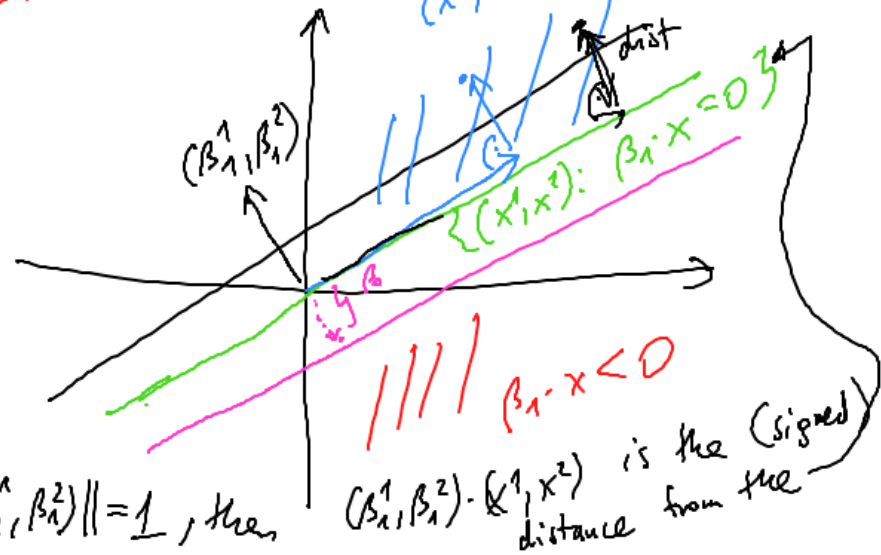
E.g. $p=2$

$$\beta_0 + \beta_1^1 X_i^1 + \beta_1^2 \cdot X_i^2 = 0$$

a ~~line~~ line in \mathbb{R}^2

$$\beta_0 = 0 \quad \beta_1 \neq 0$$

$$\underbrace{(\beta_1^1, \beta_1^2)}_{-11} \cdot (X_i^1, X_i^2) = 0 = -\beta_0$$



Then we may choose different hyperplanes that separate the classes

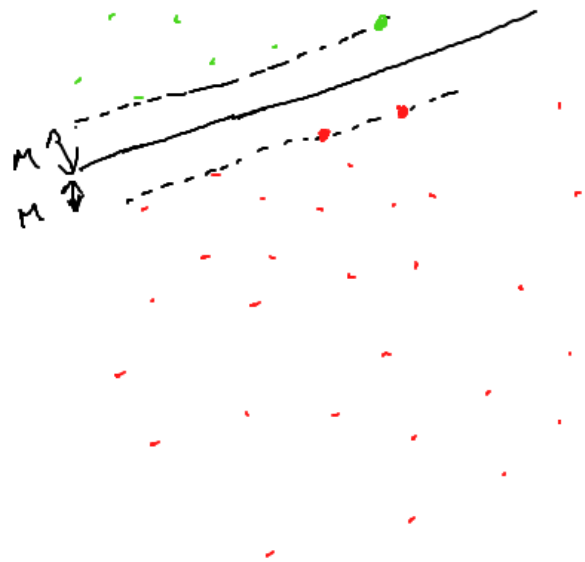
For the maximal margin classifier, we choose β_0, β_1 so that the number M (=margin) is maximised, where M is such that

$$y_i (\beta_0 + \beta_1 x_i) \geq M \quad \text{for all } i=1, \dots, n,$$

and $\|\beta_1\|^2 = 1$.

Note: just a few samples influence the final form of the classifier

What if the separating hyperplane does not exist?



C controls the bias-variance tradeoff
larger C : smaller variance (more samples influence the form of the classifier)

- How to solve this optimisation problem?

Use Karush-Kuhn-Tucker theorem (generalisation of Lagrange multipliers method)



f
 $g = \text{const.}$
 $f + \lambda g$