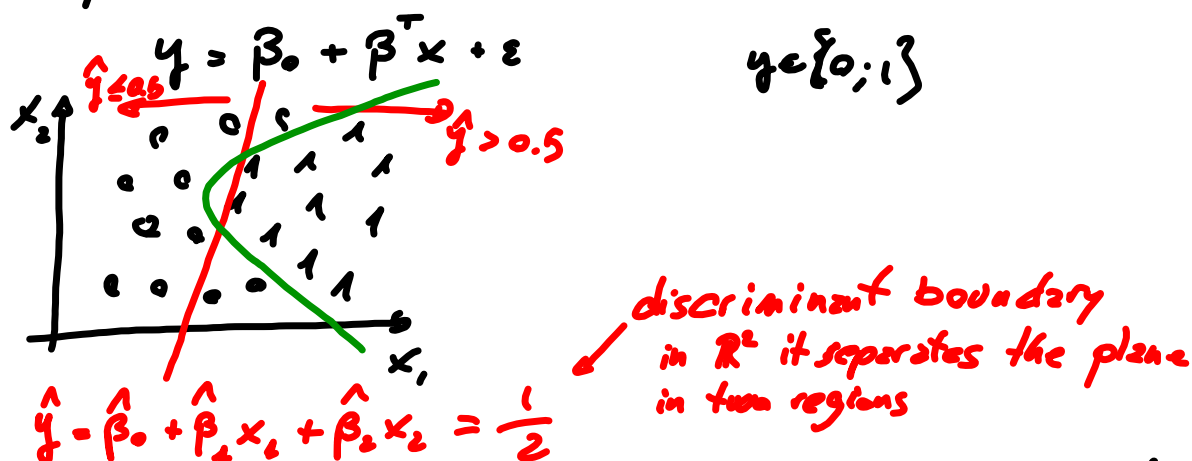


Classification via linear regression

- in logistic regression we modelled $\hat{\pi}_k(x) = \text{Pr}[Y=k|X=x]$
↳ 2 transformation
- in binary case, $y \in \{0,1\}$ $\hat{y} = 1 \Leftrightarrow \hat{\pi}_1(x) > \hat{\pi}_0(x)$
- why do not model directly y ?



We can extend this procedure by adding non-linear functions of x (e.g., polynomials, quadratic terms)

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_1^2 + \hat{\beta}_3 x_2 + \hat{\beta}_4 x_2^2 = \frac{1}{2}$$

Note:

- we often get good results, BUT
 - it is not natural to use linear Gaussian regression for classification ($y \in \{0,1\}$, $y \notin \mathbb{R}$, $\varepsilon \notin \mathcal{N}(0, \sigma^2)$)
 - all inferential tools do not work (e.g., standard errors)
 - no homoskedasticity
 - $E[\varepsilon] \neq 0$, even if we add the intercept β_0
 - masking effect

Case with several categories

- codify each class $0, \dots, K-1$ with an indicator function $Y \in \{0, \dots, K-1\}$

E.g. $K=3$

$$Y = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & 1 \end{pmatrix}$$

1st obs is in category 0

2nd and 3rd obs in c. 1

4th - 5th " " " " 2

classical multivariate linear regression

$$Y = XB + E$$

so that

$$\hat{Y} = X \underbrace{(X^T X)^{-1} X^T}_{B} Y$$

X is a $n \times p$ design matrix

B is a $p \times k$ matrix of coefficients

new observation x_0^T

$$\hat{Y}_0 = (\hat{Y}_0, \hat{Y}_1, \hat{Y}_2)$$

is assigned to the class of larger \hat{Y}_k

Note:

- we have that $\sum_{k=0}^{K-1} \hat{Y}_k = 1$, but there is no guarantee that $\hat{Y}_k, k=0, \dots, K-1 \in [0; 1]$

\Downarrow

\hat{Y}_k is not a good estimate of $\pi_k(x)$

Important issue: masking effect

(see figures 4.2 and 4.3 on pages 65 and 66 of ESL)

Possible solution: use in case of $K=3$ a quadratic term
in general: use polynomials of degree $K-1$

Discriminant Analysis

Define: Y the categorical variable which tell us to what category an observation belong

X a p -dimensional random variable

Goal: find the probability $\Pr[Y=k | X=x]$ ↖ $\hat{\pi}_k(x)$

Suppose that the population is divided in K subpopulations each of them with conditional density

$$p_k(x) = \Pr[X=x | Y=k], \quad k=0, \dots, K-1$$

Denote with π_k the probability of belonging to the subpop. k

$$\hat{\pi}_k = \Pr[Y=k] \quad \text{a priori probability of belonging to class } k$$

Then the marginal density of the whole population is

$$\begin{aligned} p(x) = \Pr[X=x] &= \sum_{k=0}^{K-1} \Pr[X=x | Y=k] \Pr[Y=k] \\ &= \sum_{k=0}^{K-1} p_k(x) \pi_k \end{aligned}$$

Then, using the Bayes theorem

$$\begin{aligned} \Pr[Y=k | X=x] &= \frac{\Pr[X=x | Y=k] \Pr[Y=k]}{\Pr[X=x]} \\ &= \frac{p_k(x) \pi_k}{\sum_{m=0}^{K-1} p_m(x) \pi_m} \end{aligned}$$

To estimate $\Pr[Y=k | X=x]$ we only need to estimate $p_k(x)$ and $\pi_k(x)$, $k=0, \dots, K-1$

- $\hat{\pi}_k(x)$ is the prior probability, so it is natural to estimate it as $\frac{n_k}{n}$ n_k is # obs belonging to class k
 n is total # obs.

Instead, we have several options to estimate $P_k(x)$

- parametric (LDA, QDA)
- non-parametric (KNN)

Since we want to assign a new observation x to the class with the highest posterior probability, it is useful to consider $\Pr[Y=k|X=x]$

$$\log \frac{\Pr[Y=k|X=x]}{\Pr[Y=m|X=x]} = \log \frac{P_k(x) \pi_k}{\sum_{c=0}^{K-1} P_c(x) \pi_c} - \log \frac{P_m(x) \pi_m}{\sum_{c=0}^{K-1} P_c(x) \pi_c}$$

$$= \log \frac{\pi_k}{\pi_m} + \log \frac{P_k(x)}{P_m(x)}$$

and the discriminant function

$$d_k(x_0) = \log \hat{\pi}_k + \log P_k(x_0)$$

we will assign x_0 to the class $\operatorname{argmax}_k d_k(x_0)$

Linear Discriminant Analysis (LDA) and Quadratic Discriminant Analysis (QDA)

LDA and QDA use multivariate Gaussian distribution for $P_k(x)$

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

$C_k \stackrel{!}{=} C$

In particular, LDA assumes $\Sigma_k = \Sigma \quad \forall k=0, \dots, K-1$

As a consequence

$$\begin{aligned} d_k(x) &= \log \hat{\pi}_k - \frac{1}{2} (x - \mu_k)^T \Sigma^{-1} (x - \mu_k) && \text{it is a scalar, so } = \mu_k^T \Sigma^{-1} x \\ &= \log \hat{\pi}_k - \frac{1}{2} (x^T \Sigma^{-1} x + \mu_k^T \Sigma^{-1} \mu_k - 2x^T \Sigma^{-1} \mu_k) \\ &= \log \hat{\pi}_k - \frac{1}{2} \mu_k^T \Sigma^{-1} \mu_k + \underline{x^T \Sigma^{-1} \mu_k} \end{aligned}$$

↳ does not depend on k

which is linear in x \rightarrow LDA

Compare two classes in terms of log-ratios

$$\begin{aligned} \log \frac{P_r[Y=k | X=x]}{P_r[Y=m | X=x]} &= \log \frac{\cancel{\exp\{-\frac{1}{2}(x-\mu_k)^T \Sigma^{-1}(x-\mu_k)\}} \hat{\pi}_k / \cancel{p(x)}}{\cancel{\exp\{-\frac{1}{2}(x-\mu_m)^T \Sigma^{-1}(x-\mu_m)\}} \hat{\pi}_m / \cancel{p(x)}} \\ &= \log \frac{\hat{\pi}_k}{\hat{\pi}_m} - \frac{1}{2} \left(\cancel{x^T \Sigma^{-1} x} - \cancel{2x^T \Sigma^{-1} \mu_k} + \mu_k^T \Sigma^{-1} \mu_k - \cancel{x^T \Sigma^{-1} x} + \cancel{2x^T \Sigma^{-1} \mu_m} + \mu_m^T \Sigma^{-1} \mu_m \right) \\ &= \log \frac{\hat{\pi}_k}{\hat{\pi}_m} - \frac{1}{2} (\mu_k - \mu_m)^T \Sigma^{-1} (\mu_k - \mu_m) + 2x^T \Sigma^{-1} (\mu_k - \mu_m) \end{aligned}$$

To compute all these quantities ($d_k(x)$, v_k), we need to plug-in the estimates of $\hat{\pi}_k$, $\hat{\mu}_k$ and $\hat{\Sigma}$

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

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

$$\hat{\Sigma} = \frac{1}{n-K} \sum_{k=0}^{K-1} \sum_{i:i_k=k} (x_i - \hat{\mu}_k)(x_i - \hat{\mu}_k)^T$$

** parameters to estimate $p_k + \frac{p(p+1)}{2}$
 μ_k is p -dimensional*

To have curved boundaries, we should add quadratic terms

$$x_i = (x_{i1}, x_{i2}) \longrightarrow x_i = (x_{i1}, x_{i2}, x_{i1}^2, x_{i2}^2, x_{i1}x_{i2})$$

- with $K=2$, there is no big differences in using LDA or linear regression
- with $K > 2$, substantial differences (LDA does not suffer from the masking effect issue)

QDA \rightarrow we remove the condition $\Sigma_k = \Sigma \forall k$

$$d_k(x) = \log \hat{\pi}_k - \frac{1}{2} \underbrace{(x - \mu_k)^T \Sigma_k^{-1} (x - \mu_k)}_{\substack{\text{we cannot remove the quadratic term} \\ x^T \Sigma_k^{-1} x}} - \frac{1}{2} \log (\det(\Sigma_k))$$

\downarrow
QLA

To estimate $d_k(x)$, the same estimates for $\hat{\pi}_k$ and $\hat{\mu}_k$, but

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

$\#$ parameters to estimate
 $pK + \frac{p(p+1)}{2} K$

Note:

- as often, the simpler model (LDA) performs better than the more complex (QDA)
- LDA needs the estimation of fewer parameters (better use of the information)