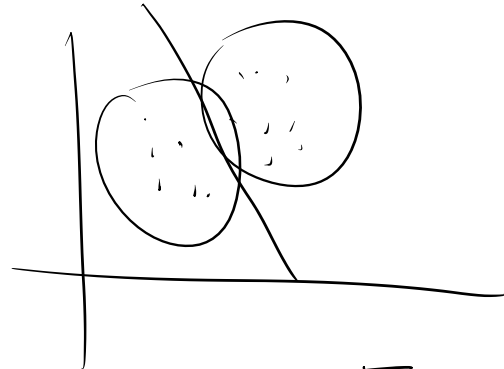
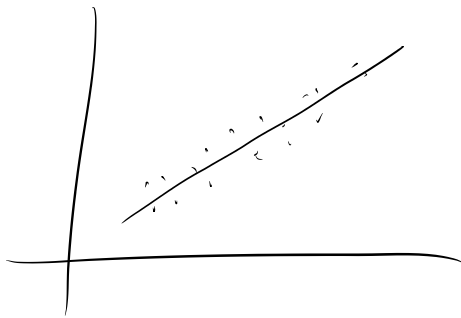


		Clasif.	
		π_1	π_2
Verdader	π_1	0	$C(2 1)$
	π_2	$C(1 2)$	0

El costo esperado o promedio de clasif

$$ECM = C(2|1) \cdot p(2|1) p_1 + C(1|2) p(1|2) p_2 \rightarrow \text{pepe}$$

Buscamos minimizar ese costo.



$p(2|1)$ prob. de clasif. un objeto en π_2 siendo π_1

$p(1|2)$ " " " " " " " π_1 siendo π_2

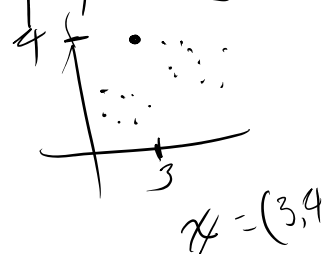
Nuestro objetivo es minimizar a pepe

Teorema:

Las regiones R_1 y R_2 que minimizan a pepe se definen por los x que satisfacen

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \left(\frac{C(1|2)}{C(2|1)} \right) \cdot \left(\frac{p_2}{p_1} \right)$$

$$R_2: \frac{f_1(x)}{f_2(x)} < \left(\frac{C(1|2)}{C(2|1)} \right) \cdot \left(\frac{p_2}{p_1} \right)$$



$$f_1(x) = 3x_1 + 2x_2$$

$$f_2(x) \quad | \quad C(2|1) \quad | \quad p_1$$

Ejercicio: 11.3 libro. (Tarea)

Para usar esta técnica necesitamos:

- 1) Las funciones de densidad (por lo menos evaluadas en los puntos de interés x)
- 2) Los costos de clasificación errónea
- 3) Las probabilidades previas.

Caso especial,

$$\text{si } \frac{p_2}{p_1} = 1 \quad (p_2 = p_1)$$

$$R_1: \frac{f_1(x)}{f_2(x)} \geq \frac{c(1|2)}{c(2|1)}$$

$$R_2: \frac{f_1(x)}{f_2(x)} < \frac{c(1|2)}{c(2|1)}$$

$$2) c(1|2) = c(2|1) \Rightarrow \frac{c(1|2)}{c(2|1)} = 1$$

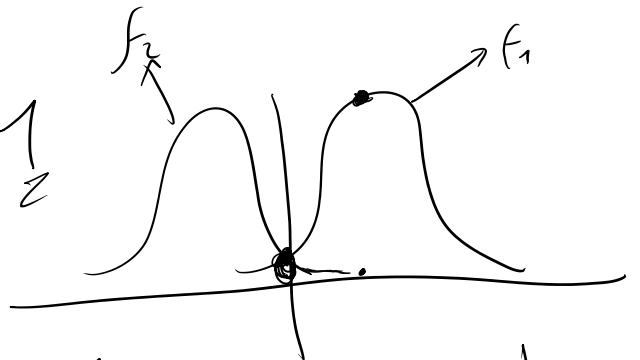
$$R_1: \frac{f_1(x)}{f_2(x)} \geq \frac{p_2}{p_1}$$

$$R_2: \frac{f_1(x)}{f_2(x)} < \frac{p_2}{p_1}$$

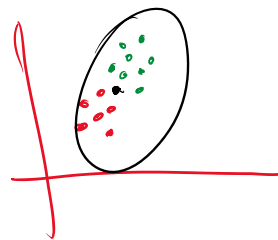
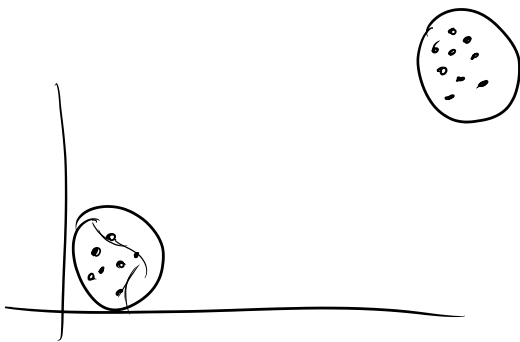
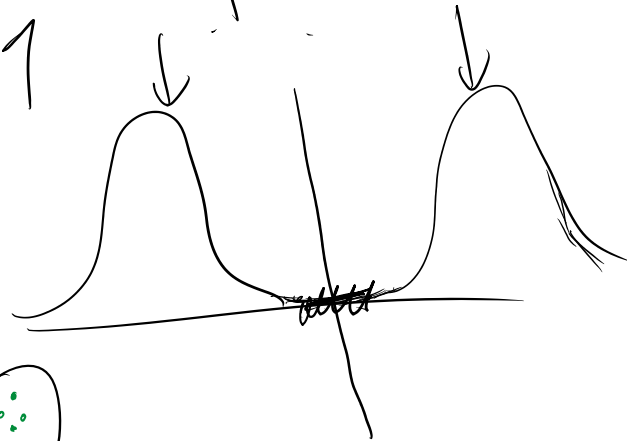
$$f_2(x) \quad p_1$$

3) costos iguales y p_1, p_2 también

$$R_1: \frac{f_1(x)}{f_2(x)} \geq 1$$



$$R_2: \frac{f_1(x)}{f_2(x)} < 1$$



Clasificar en 2 pobl. normales multivariadas

Bajo la supo. de normalidad la clasif. es mucho más sencilla y resulta más práctica

Supongamos $f_1(x), f_2(x)$ PDF normales multiv.
con medias μ_1 y μ_2 , cov Σ_1, Σ_2 respect.

Caso 1 $\Sigma_1 = \Sigma_2 = \Sigma$

$$f_i(x) = \frac{1}{(2\pi)^{p/2} |\Sigma|^{1/2}} \cdot \exp\left(-\frac{1}{2} (x - \mu_i)' \Sigma^{-1} (x - \mu_i)\right)$$

Supongamos μ_1, μ_2 y Σ conocidos

Teorema

$$R_1: \exp\left(-\frac{1}{2} (X-\mu_1)' \Sigma^{-1} (X-\mu_1) + \frac{1}{2} (X-\mu_2)' \cdot \Sigma^{-1} (X-\mu_2)\right) \geq \frac{C(1|2)}{C(2|1)} \cdot \frac{P_2}{P_1}$$

$$R_2: \boxed{} < \frac{C(1|2)}{C(2|1)} \frac{P_2}{P_1}$$

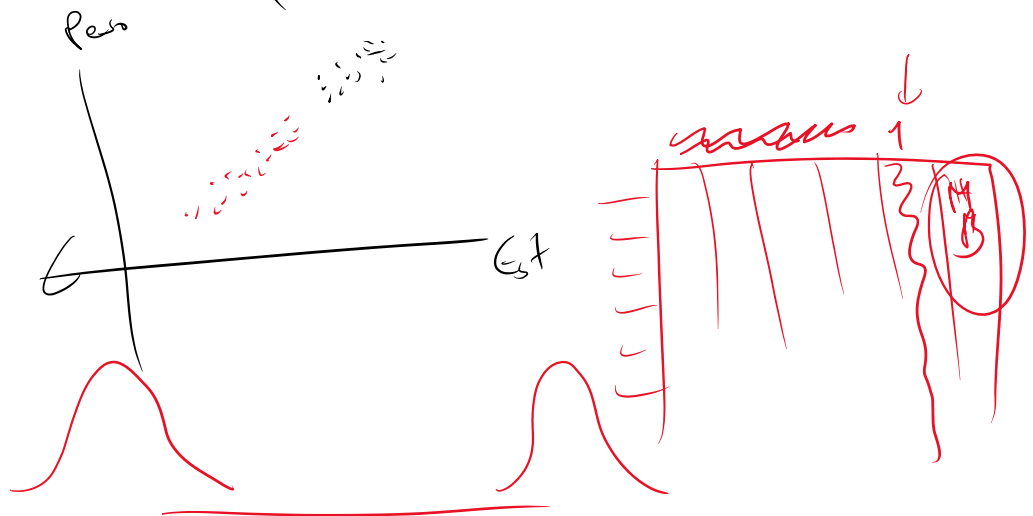
Teorema:

Sean π_1, π_2 pobl. con PDF normal multiv., entonces la regla de clasif. que minimiza el error es:

Clasificar X_0 como π_1 si:

$$(\mu_1 - \mu_2)' \Sigma^{-1} X_0 - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln \left(\frac{C(1|2)}{C(2|1)} \cdot \frac{P_2}{P_1} \right)$$

Demstración (libro - TARE A)



En la mayoría de los casos μ_1, μ_2 y Σ son

En la mayoría de los casos μ_1 , μ_2 y Σ son desconocidos y debemos estimarlos con los posibles errores que eso implica.

Supongamos que hay n_1 obs. de $X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ de π_1
 y que hay n_2 obs. de $X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}$ de π_2
 # de variables

con $n_1 + n_2 - 2 \geq \hat{p}$.

tenemos $X_1 = \begin{pmatrix} x_{11}' \\ \vdots \\ x_{1n_1}' \end{pmatrix}$, $X_2 = \begin{pmatrix} x_{21}' \\ \vdots \\ x_{2n_2}' \end{pmatrix}$

con medias \bar{X}_1 y \bar{X}_2

y cov. muestrales

$$S_1 = \frac{1}{n_1 - 1} \sum_{i=1}^{n_1} (x_{1i} - \bar{X}_1)(x_{1i} - \bar{X}_1)'$$

S_2 : igual pero con X_2

$$S_{\text{pooled}} = \left(\frac{n_1 - 1}{n_1 + n_2 - 2} \right) S_1 + \left(\frac{n_2 - 1}{n_1 + n_2 - 2} \right) S_2$$

↓
 estimador insesgado de Σ

Entonces

Clasificar X_0 como π_1 si:

$$\left((\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} X_0 - \frac{1}{2} (\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} (\bar{X}_1 + \bar{X}_2) \right) \geq \ln \left(\frac{P_1}{P_2} \right)$$

$$\boxed{\text{Jose}} \geq \frac{1}{2} \ln \left(\frac{C(1|2)}{C(2|1)} \right) \frac{P_2}{P_1}$$

$$\text{Si } \frac{C(1|2)}{C(2|1)} \frac{P_2}{P_1} = 1 \Rightarrow \ln(1) = 0$$

Entonces, clasificamos como Π_1 si $\boxed{\text{Jose}} \geq 0$

$$\begin{aligned} \text{Sea } \hat{y}_0 &= (\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} X_0 \\ &= \hat{\alpha}' X_0 \end{aligned}$$

$$\text{y también } \hat{m} = \frac{1}{2} (\bar{X}_1 - \bar{X}_2)' S_{\text{pooled}}^{-1} (\bar{X}_1 + \bar{X}_2) = \frac{1}{2} \underbrace{\left(\begin{matrix} \hat{\alpha}' \bar{X}_1 \\ \hat{\alpha}' \bar{X}_2 \end{matrix} \right)}_{\hat{\alpha}' \bar{X}_0} (\bar{X}_1 + \bar{X}_2)$$

Usaremos \hat{y} y \hat{m} para clasificar

(Si \hat{y}_0 es mayor o menor al promed. de \bar{y}_1 y \bar{y}_2)

costs) This example is adapted from a study [4] concerned with the detection of hemophilia A carriers. (See also Exercise 11.32.)

To construct a procedure for detecting potential hemophilia A carriers, blood samples were assayed for two groups of women and measurements on the two variables,

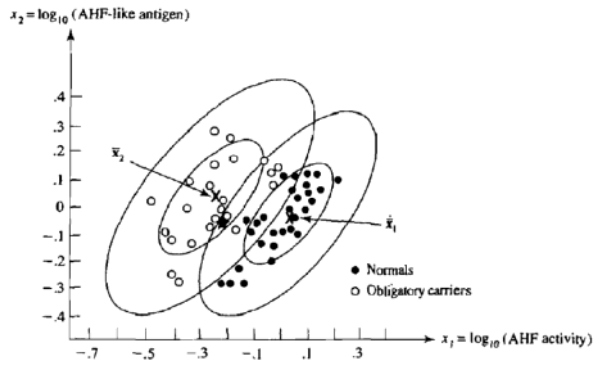
$$X_1 = \log_{10}(\text{AHF activity})$$

$$X_2 = \log_{10}(\text{AHF-like antigen})$$

recorded. ("AHF" denotes antihemophilic factor.) The first group of $n_1 = 30$ women were selected from a population of women who did not carry the hemophilia gene. This group was called the *normal* group. The second group of $n_2 = 22$ women was selected from known hemophilia A carriers (daughters of hemophiliacs, mothers with more than one hemophilic son, and mothers with one hemophilic son and other hemophilic relatives). This group was called the *obligatory carriers*. The pairs of observations (x_1, x_2) for the two groups are plotted in Figure 11.4. Also shown are estimated contours containing 50% and 95% of the probability for bivariate normal distributions centered at \bar{x}_1 and \bar{x}_2 , respectively. Their common covariance matrix was taken as the pooled sample covariance matrix S_{pooled} . In this example, bivariate normal distributions seem to fit the data fairly well.

The investigators (see [4]) provide the information

$$\bar{x}_1 = \begin{bmatrix} -.0065 \\ -.0390 \end{bmatrix}, \quad \bar{x}_2 = \begin{bmatrix} -.2483 \\ .0262 \end{bmatrix}$$



and

$$S_{\text{pooled}}^{-1} = \begin{bmatrix} 131.158 & -90.423 \\ -90.423 & 108.147 \end{bmatrix}$$

Therefore, the equal costs and equal priors discriminant function [see (11-19)] is

$$\begin{aligned} \hat{y} &= \hat{a}'x = [\bar{x}_1 - \bar{x}_2]'S_{\text{pooled}}^{-1}x \\ &= \begin{bmatrix} .2418 & -.0652 \end{bmatrix} \begin{bmatrix} 131.158 & -90.423 \\ -90.423 & 108.147 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \\ &= 37.61x_1 - 28.92x_2 \end{aligned}$$

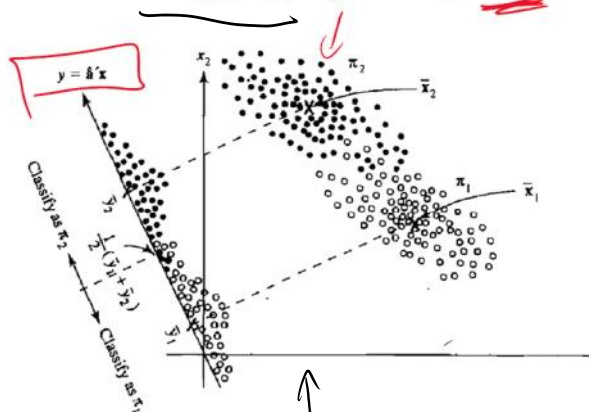
Moreover,

$$\begin{aligned} \bar{y}_1 &= \hat{a}'\bar{x}_1 = \begin{bmatrix} 37.61 & -28.92 \end{bmatrix} \begin{bmatrix} -.0065 \\ -.0390 \end{bmatrix} = .88 \\ \bar{y}_2 &= \hat{a}'\bar{x}_2 = \begin{bmatrix} 37.61 & -28.92 \end{bmatrix} \begin{bmatrix} -.2483 \\ .0262 \end{bmatrix} = -10.10 \end{aligned}$$

and the midpoint between these means [see (11-20)] is

and the midpoint between these means [see (11-20)] is

$$\hat{m} = \frac{1}{2}(\bar{y}_1 + \bar{y}_2) = \frac{1}{2}(.88 - 10.10) = -4.61$$



Método de discriminante de Fisher.