

## Clasificación 2

Sunday, May 8, 2022 6:39 PM

## ## Clasificación 2 ##

Clasificación con 2 poblaciones multivariadas

Bajo la suposición de normalidad, la clasificación es sencilla y resulta práctica (incluso si son aprox normales).  
 Suponemos  $f_1(x)$ ,  $f_2(x)$  PDF's normal multivar con medias  $\mu_1$  y  $\mu_2$ , y cov.  $\Sigma_1$ ,  $\Sigma_2$ .

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

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

Sup. q'  $\mu_1$ ,  $\mu_2$  y  $\Sigma$  conocidos

Teo:

$$R_1 = \exp\left(-\frac{1}{2} (x - \mu_1)' \Sigma^{-1} (x - \mu_1) + \frac{1}{2} (x - \mu_2)' \Sigma^{-1} (x - \mu_2)\right) \\ \geq \frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1} \quad \leftarrow \frac{f(x_1)}{f(x_2)}$$

$$R_2 = \frac{f(x_1)}{f(x_2)} < \frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1}$$

Teo: Sean  $\pi_1$  y  $\pi_2$  poblaciones con PDF normal multivar. entonces la regla de clasificación q' minimiza el error es: $x_0$  se clasifica como  $\pi_1$  si:

$$(\mu_1 - \mu_2)' \Sigma^{-1} x_0 - \frac{1}{2} (\mu_1 - \mu_2)' \Sigma^{-1} (\mu_1 + \mu_2) \geq \ln\left(\frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1}\right)$$

Obs: Esto sirve si conocemos  $\mu_1$ ,  $\mu_2$  y  $\Sigma$ , pero en muchos casos no lo conocemos.

Sup. q' hay  $n_1$  obs de  $x = \begin{bmatrix} x_1 \\ x_p \end{bmatrix}$  de  $\pi_1$  con  $n_1 + n_2 - 2 \geq p$   
 $n_2$  obs. de  $x = \begin{bmatrix} x_1 \\ x_p \end{bmatrix}$  de  $\pi_2$

$x_1 = \begin{bmatrix} x_{11} \\ \vdots \\ x_{1n_1} \end{bmatrix}$   $x_2 = \begin{bmatrix} x_{21} \\ \vdots \\ x_{2n_2} \end{bmatrix}$  con medias muestrales  $\bar{x}_1$ ,  $\bar{x}_2$   
 y cov. muestrales:

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

Est. interrogado de  $x$   
(seguimos asumiendo q' es en  $\pi_1$  y  $\pi_2$ )

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

$$S_2 = \frac{1}{n_2 - 1} \sum_{j=1}^{n_2} (x_{2j} - \bar{x}_2)(x_{2j} - \bar{x}_2)'$$

Clasificar  $x_0$  como  $\pi_1$  si:

$$(\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} x_0 - \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} (\bar{x}_1 + \bar{x}_2) \geq \ln\left(\frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1}\right)$$

Clasificar  $x_0$  como  $\pi_1$  si no.Obs: Si  $\frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1} = 1$   $\ln(1) = 0$ .

$$\text{sea } \hat{y}_0 = (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} x_0 \\ = \hat{a}' x_0$$

$$\hat{m} = \frac{1}{2} (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} (\bar{x}_1 + \bar{x}_2)$$

$$= \frac{1}{2} (\bar{y}_1 + \bar{y}_2) \quad \text{donde } \bar{y}_1 = \hat{a}' \bar{x}_1 \\ \bar{y}_2 = \hat{a}' \bar{x}_2$$

Podemos clasificar comparando  $\hat{y}_0$  con  $\hat{m}$   
 si es mayor o menor del promedio de medias  $\bar{y}_1$ ,  $\bar{y}_2$   
 si es  $\geq \hat{m}$ , es  $\pi_1$

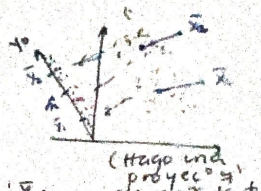
Regla de clasificación usando la discriminante de Fisher

Clasifica  $x_0$  como  $\pi_1$  si:

$$\hat{y}_0 = (\bar{x}_1 - \bar{x}_2)' S_{pooled}^{-1} x_0 \geq \hat{m} \quad \text{o} \quad \hat{y}_0 - \hat{m} \geq 0$$

Si  $\hat{y}_0 - \hat{m} < 0$ , se clasifica como  $\pi_2$ Obs:  $\hat{y}$  (min error cuando los costos y la prob previas = 1) es la  $f_0$  de Fisher q' maximiza la separación.Caso 2:  $\Sigma_1 \neq \Sigma_2$ 

$$R_1 = -\frac{1}{2} x' (\Sigma_1^{-1} - \Sigma_2^{-1}) x + (\mu_1' \Sigma_1^{-1} - \mu_2' \Sigma_2^{-1}) x - k \geq \ln\left(\frac{C(1|2)}{C(2|1)} \cdot \frac{p_2}{p_1}\right)$$



... es menor,  $\mu_2$ .

$$K = \frac{1}{2} \ln \left( \frac{\det(\Sigma_1)}{\det(\Sigma_2)} \right) + \frac{1}{2} (\mu_1' \Sigma_1^{-1} \mu_1 - \mu_2' \Sigma_2^{-1} \mu_2)$$

Como en la practica no se conocen, estimamos  $\mu_1, \mu_2, \Sigma_1$  y  $\Sigma_2$ .

Se clasifica  $x_0$  como  $\pi_1$  ( $\pi_2$ ) si:

$$-\frac{1}{2} x_0' (S_1^{-1} - S_2^{-1}) x_0 + (\bar{x}_1' S_1^{-1} - \bar{x}_2' S_2^{-1}) x_0 - K \geq \ln \left( \frac{C(1|2)}{C(2|1)} \right) \frac{P_2}{P_1}$$

$S_1 < S_2$ .

**evaluación de clasificación**

Podemos evaluar la calidad de la clasificación con tasas de error (prueba de clasificación incorrecta).

i) Si se conoce la  $\sim$  de la población, es directo.

ii) Si no, usar regla de clasificación muestral y se evalúa con muestras.

**Para i)**  $TPM = P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$  (Total Probable Mistake).  
Es el OER (Optimal Error Rate).

**Para ii)** Algunos parámetros desconocidos.

**Caso 1:** Se conocen  $f_1, f_2$ .

$$\text{AER (Actual error rate)} = P_1 \int_{R_2} f_1(x) dx + P_2 \int_{R_1} f_2(x) dx$$

$R_2 \rightarrow \text{def. con } n_1 \text{ y } n_2.$

**Caso 2:**  $f_1, f_2$  desconocidas.

APER (apparent error rate)  $\rightarrow$  matriz de confusión: pop. real

	Pobl 1	Pobl 2
$\pi_1$	$n_{11}$	$n_{12}$
$\pi_2$	$n_{21}$	$n_{22}$

$n_{11} + n_{12} = n_1$   
 $n_{21} + n_{22} = n_2$

$$\text{APER} = \frac{n_{12} + n_{21}}{n_1 + n_2} : (\text{entre menor el APER, mejor } \omega).$$

Se necesitan  $n_1$  y  $n_2$  muy grandes.

2 grupos  $\rightarrow$  entrenamiento.  
 $\rightarrow$  validación.

