

### 3. ANÁLISIS DE CORRESPONDENCIA

#### • PREGUNTA 1

(a) Plantea los problemas de independencia de caracteres y homogeneidad de poblaciones.

##### \* INDEPENDENCIA DE CARACTERES

Consideramos dos variables categóricas  $(A, B)$ , donde

$$\left. \begin{array}{l} A \text{ con } n \text{ modalidades : } A_1, \dots, A_n \\ B \text{ con } p \text{ modalidades : } B_1, \dots, B_p \end{array} \right\}$$

El problema consiste en comprobar la hipótesis de independencia  
 $H_0 : P(A_i \cap B_j) = P(A_i) P(B_j), \forall i, j$

(Para ello, seleccionamos una muestra  $m$  de tamaño  $N$  y observamos  $(A, B)$  en cada elemento).

##### \* HOMOGENEIDAD DE POBLACIONES

Consideramos una variable categórica  $B$  con  $p$  modalidades  $B_1, \dots, B_p$ , y estudiamos la homogeneidad de  $B$  en  $n$  poblaciones  $(P_1, \dots, P_n)$ .

El problema consiste en comprobar la hipótesis de homogeneidad  
 $H_0 : P_1(B_j) = \dots = P_n(B_j) (= P(B_j)), j = 1, \dots, p.$

(Para ello seleccionamos de forma independiente una muestra  $m_i$  de tamaño  $N_i$  de cada población).

(b) Demuestra que los estadísticos  $\chi^2$  asociados a ambos problemas y sus distribuciones coinciden.

Para resolver el problema de independencia, seleccionamos una muestra  $m$  de tamaño  $N$  y observamos  $(A, B)$  en cada elemento.

$$N_{ij} \equiv \text{nº de elementos de } m \text{ que presentan } A_i \cap B_j \\ (\sim B_i(N, P(A_i \cap B_j)))$$

donde, la tabla de contingencia (matriz de datos experimentales) es:

$$(*) \quad \begin{array}{c} \begin{array}{cccc} & B_1 & \dots & B_j & \dots & B_p & \text{Marginal de A} \\ A_1 & N_{11} & \dots & N_{1j} & \dots & N_{1p} & N_{1.} \\ \vdots & \vdots & & \vdots & & \vdots & \vdots \\ A_i & N_{i1} & \dots & N_{ij} & \dots & N_{ip} & N_{i.} = \sum_{j=1}^p N_{ij} \\ \vdots & \vdots & & \vdots & & \vdots & \vdots \\ A_n & N_{n1} & \dots & N_{nj} & \dots & N_{np} & N_{n.} \end{array} \\ \text{Marginal de B} \quad N_{.1} \dots N_{.j} = \sum_{i=1}^n N_{ij} \dots N_{.p} \end{array}$$

Luego, el estadístico chi-cuadrado para la hipótesis de independencia es:

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^p \frac{(N_{ij} - e_{ij})^2}{e_{ij}}, \quad \text{con } e_{ij} = \hat{E}_{H_0}(N_{ij})$$

Como  $N_{ij} \underset{H_0}{\sim} Bi(N, P(A_i)P(B_j))$ , entonces:

$$e_{ij} = N \hat{P}_{H_0}(A_i) \hat{P}_{H_0}(B_j) = N \cdot \frac{N_{i.}}{N} \cdot \frac{N_{.j}}{N} = \frac{N_{i.} \cdot N_{.j}}{N}$$

$$\text{Luego, } \chi^2 \xrightarrow{L} \chi^2_{(n-1)(p-1)} \quad \begin{array}{l} [n(p-1) - (p-1)] \\ \text{var. lin indep - lin indep bajo } H_0 \end{array}$$

Para resolver el pb de homogeneidad seleccionamos de forma indep. una muestra  $m_i$  de tamaño  $N_i$  de cada población. Ahora,

$N_{ij} \equiv$  n° de elementos de  $m_i$  que presentan  $B_j$  ( $\sim Bi(N_i, P_i(B_j))$ )

$$\begin{array}{c} \begin{array}{cccc} & B_1 & \dots & B_j & \dots & B_p \\ P_1 & N_{11} & \dots & N_{1j} & \dots & N_{1p} \\ \vdots & \vdots & & \vdots & & \vdots \\ P_i & N_{i1} & \dots & N_{ij} & \dots & N_{ip} \\ \vdots & \vdots & & \vdots & & \vdots \\ P_n & N_{n1} & \dots & N_{nj} & \dots & N_{np} \end{array} \\ N_{.1} \dots N_{.j} \dots N_{.p} \end{array} \quad \begin{array}{l} N_{1.} = N_1 \\ \vdots \\ N_{i.} = N_i \\ \vdots \\ N_{n.} = N_n \end{array} \quad N_{.j} = \sum_{i=1}^n N_{ij}$$

El estadístico chi-cuadrado para la hipótesis de homogeneidad es:

$$\chi^2 = \sum_{i=1}^n \sum_{j=1}^p \frac{(N_{ij} - e_{ij})^2}{e_{ij}}, \quad \text{con } e_{ij} = \hat{E}_{H_0}(N_{ij})$$

Como  $N_{ij} \underset{H_0}{\sim} B_i(N_i, P_i(B_j)) \Rightarrow e_{ij} = N_i \cdot \hat{P}_{H_0}(B_j) = N_i \cdot \frac{N_{.j}}{N}$

Luego,  $\chi^2 \xrightarrow{\mathcal{L}} \chi^2_{(n-1)(p-1)} \leftarrow [n(p-1)] - (p-1)$

## • PREGUNTA 2

Sean  $r'_1, \dots, r'_n$  las distribuciones condicionadas por filas asociadas a una tabla de contingencia.

(\*)

(a) Define el centro de gravedad de las filas:  $m'_r$

Transformamos la tabla de contingencia para obtener las distribuciones condicionadas por filas:

$$M_r = \begin{bmatrix} r'_1 \\ \vdots \\ r'_i \\ \vdots \\ r'_n \end{bmatrix} = \begin{bmatrix} \frac{N_{11}}{N_{1.}} & \dots & \frac{N_{1j}}{N_{1.}} & \dots & \frac{N_{1p}}{N_{1.}} \\ \vdots & & \vdots & & \vdots \\ \frac{N_{i1}}{N_{i.}} & \dots & \frac{N_{ij}}{N_{i.}} & \dots & \frac{N_{ip}}{N_{i.}} \\ \vdots & & \vdots & & \vdots \\ \frac{N_{n1}}{N_{n.}} & \dots & \frac{N_{nj}}{N_{n.}} & \dots & \frac{N_{np}}{N_{n.}} \end{bmatrix} \begin{bmatrix} \frac{N_{1.}}{N} \\ \vdots \\ \frac{N_{i.}}{N} \\ \vdots \\ \frac{N_{n.}}{N} \end{bmatrix}$$

$$m'_r = \begin{bmatrix} \frac{N_{.1}}{N} & \dots & \frac{N_{.j}}{N} & \dots & \frac{N_{.p}}{N} \end{bmatrix}$$

Llamamos  $r_{ij} = \frac{N_{ij}}{N_{i.}}$ ,  $f_{i.} = \frac{N_{i.}}{N}$ ,  $f_{.j} = \frac{N_{.j}}{N}$

Por tanto, definimos el centro de gravedad  $m_r$  de las filas de  $M_r$  como:

$$m'_r = \sum_{i=1}^n f_{i.} r'_i = \begin{bmatrix} f_{.1} & \dots & f_{.j} & \dots & f_{.p} \end{bmatrix}$$

(b) Define la distancia chi-cuadrado entre dos distribuciones.

La distancia chi-cuadrado es una distancia entre distribuciones de probabilidad.

(Distancia chi-cuadrado entre dos filas de  $M_r$ ):

$$d_{\chi^2}(r_i, r_{i'}) = \sum_{j=1}^p \frac{(r_{ij} - r_{i'j})^2}{f_{.j}} = (r_i - r_{i'})' D_p^{-1} (r_i - r_{i'})$$

donde,  $D_p = \text{diag}(f_{.1}, \dots, f_{.p})$

(c) Demuestra que  $\chi^2 = N \sum_{i=1}^n f_{i.} d_{\chi^2}(r_i, m_r)$

$$\begin{aligned}
 \chi^2 &= \sum_{i=1}^n \sum_{j=1}^p \frac{(N_{ij} - e_{ij})^2}{e_{ij}} = \sum_{i=1}^n \sum_{j=1}^p \frac{\left( N_{ij} - \frac{N_{i.} N_{.j}}{N} \right)^2}{\frac{N_{i.} N_{.j}}{N}} = \\
 &= \sum_{i=1}^n \sum_{j=1}^p \frac{N_{i.}}{f_{i.} N} \left( \frac{N_{ij}}{N_{i.}} - \frac{N_{.j}}{N} \right)^2 = \\
 &\quad \text{sacamos } N_{i.} \quad \begin{array}{c} r_{ij} = \frac{N_{ij}}{N_{i.}} \\ e_{ij} \text{ bajo } H_0 = \frac{N_{.j}}{N} \end{array} \\
 &= N \sum_{i=1}^n f_{i.} \underbrace{\sum_{j=1}^p \frac{1}{f_{.j}} (r_{ij} - f_{.j})^2}_{d_{\chi^2}(r_i, m_r)} = N \sum_{i=1}^n f_{i.} d_{\chi^2}(r_i, m_r)
 \end{aligned}$$

### • PREGUNTA 3

Describe el objetivo del análisis de correspondencia (por filas).

El objetivo del análisis de correspondencia por filas es representar los perfiles filas en un espacio de menor dimensión (generalmente 2) de forma que si

$$M_r = \begin{pmatrix} r_{11} & \dots & r_{1p} \\ \vdots & \ddots & \vdots \\ r_{n1} & \dots & r_{np} \end{pmatrix} \xrightarrow{\text{transf.}} \begin{pmatrix} P_{r11} & \dots & P_{r1p} \\ \vdots & & \vdots \\ P_{rn1} & \dots & P_{rnp} \end{pmatrix} = PM_r$$

buscamos la base de  $\mathbb{R}^p$ ,  $\{u_1, \dots, u_p\}$  de esa transformación de forma que se verifique:

$$* d_{\chi^2}(r_i, r_j) = d_e^2(P_{ri}, P_{rj})$$

Es decir, los perfiles próximos en la métrica chi-cuadrado tienen representaciones próximas en la distancia euclídea.

\* Sean  $v_1, \dots, v_p$  las varianzas de las columnas de  $PM_r$ . Debe cumplirse que  $v_1 \geq v_2 \geq \dots \geq v_p$  y que las últimas varianzas sean próximas a 0.

$$(*) \text{ NOTA : } P_{ri}^t = (r_i^t D_p^{-1} u_1, \dots, r_i^t D_p^{-1} u_p)$$

#### • PREGUNTA 4

Analogías y diferencias entre el ACP y el AC.

- \* Ambos procedimientos (tanto ACP como AC) buscan reducir la dimensión del espacio original de estudio. Para ello, se realizan transformaciones (cambios de base) que mantengan la máxima varianza en las primeras componentes o coordenadas.
- \* Sin embargo, en el ACP se verifica que la distancia euclídea de los individuos originales y transformados coinciden, mientras que en el AC, la distancia chi-cuadrado de las coordenadas originales (distribuciones) coincide con el cuadrado de la distancia euclídea de las coordenadas en la nueva base.  
Además, en el ACP, las coordenadas principales deben estar incorreladas. (componentes)

### • PREGUNTA 5

Plantea y resuelve el problema de optimización al que conduce el AC.

El problema que se plantea es determinar los ejes  $u_1, \dots, u_p$  unitarios y ortogonales, que maximicen la varianza (ponderada) de las proyecciones.

(1) Primer eje ( $u_1$ )

$$\begin{aligned} \max_{u_1} \quad & u_1^t D_p^{-1} S_r u_1 \\ \text{s.a.} \quad & u_1^t D_p^{-1} u_1 = 1 \end{aligned}$$

(2) segundo eje ( $u_2$ )

$$\begin{aligned} \max_{u_2} \quad & u_2^t D_p^{-1} S_r u_2 \\ \text{s.a.} \quad & u_2^t D_p^{-1} u_2 = 1 \\ & u_2^t D_p^{-1} u_1 = 0 \end{aligned}$$

Así se procede con los demás  $u_k$ ,  $k=1, \dots, p$ .

Solución de (1):

$$\text{Sea } L(u_1, \lambda) = u_1^t D_p^{-1} S_r u_1 - \lambda (u_1^t D_p^{-1} u_1 - 1)$$

$$\frac{\partial L}{\partial u_1} = 2 D_p^{-1} S_r u_1 - 2 \lambda D_p^{-1} u_1 = 0_{(p)}$$

$$\Rightarrow D_p^{-1} S_r u_1 = D_p^{-1} \lambda u_1 \Rightarrow \boxed{S_r u_1 = \lambda u_1} \quad (\lambda, u_1) \text{ autovector / vec de } S_r$$

Multiplicando por  $u_1^t D_p^{-1}$ :

$$u_1^t D_p^{-1} S_r u_1 = u_1^t D_p^{-1} \lambda u_1 = \lambda \underbrace{u_1^t D_p^{-1} u_1}_1$$

$$\Rightarrow u_1^t D_p^{-1} S_r u_1 = \lambda$$

Como estamos maximizando,  $\lambda \equiv \text{máx. autovector de } S_r$

$\Rightarrow u_1 \equiv \text{autovector asociado al máximo autovector de } S_r$ .



## • PREGUNTA 7

Concepto de vértice en el problema del AC por filas.  
Interpretación.

Los vértices de las filas son una distribución de probabilidad degenerada que concentra toda la información en la  $i$ -ésima fila o categoría.

$$(B_1) v_{f_1}^t = (1, 0, \dots, 0)_{p \times 1} \rightarrow \text{concentra toda la probabilidad en la 1ª categoría}$$

$$\vdots$$

$$(B_i) v_{f_i}^t = (0, \dots, \underset{\downarrow i}{1}, 0, \dots, 0)_{p \times 1}$$

"Cada vértice  $v_{f_i}$  es una distrib. degenerada que concentra toda su información en  $B_i$ "

$$\vdots$$

$$(B_p) v_{f_p}^t = (0, \dots, 0, 1)_{p \times 1} \quad (\text{Los vértices son distribuciones extremas}).$$

Interpretación:

Como  $d_{\chi^2}(r_i, v_{f_i}) = d_e^2(P_{r_i}, P_{v_{f_i}})$  podemos interpretar la distancia del perfil fila  $r_i$  a la característica  $B_j$ .

### • PREGUNTA 8

Determina las coordenadas de  $m_r'$  en la nueva base.

Veamos que  $(0, m_r)$  es un autovalor / autovector de  $S_r$ .  
Para ello, tenemos que comprobar que se verifica

$$S_r \cdot m_r = 0_p$$

Entonces,

$$S_r \cdot m_r = \sum_{i=1}^n f_i \cdot (r_i - m_r)(r_i - m_r)^t D_p^{-1} m_r = \dots$$

$$\begin{pmatrix} \frac{1}{f \cdot 1} & & \\ & \ddots & \\ & & \frac{1}{f \cdot p} \end{pmatrix} \begin{pmatrix} \frac{0}{f \cdot 1} \\ \vdots \\ \frac{p}{f \cdot p} \end{pmatrix} = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}_p = \mathbb{1}_p$$

Además,

$$\left. \begin{array}{l} r_i^t \cdot \mathbb{1}_p = 1 \\ m_r^t \cdot \mathbb{1}_p = 1 \end{array} \right\} \Rightarrow \text{Por tanto, } (r_i - m_r)^t D_p^{-1} m_r = 0 \text{ (escalar)}$$

Luego,

$$\dots = 0_p.$$

$\Rightarrow$  como  $(0, m_r)$  es autovalor / autovector de  $S_r$ ,  $m_r$  es un eje de la nueva base  $\Rightarrow m_r' D_p^{-1} u_k = 0$  (para  $u_k \neq m_r$ )

$\Rightarrow (0, \dots, 0, 1)$  son las coordenadas de  $m_r^t$ .

Las coordenadas de  $r_i$  sobre  $m_r$ :

$$\begin{bmatrix} r_1^t D_p^{-1} m_r \\ \vdots \\ r_n^t D_p^{-1} m_r \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}$$

(Coordenadas de  $m_r$ )

Cuando proyectamos todos tienen coordenada 1.

Es constante ya que la varianza de dicha columna vale 0.

Para la solución de (2):

$$\text{Sea } L(u_2, \beta_1, \beta_2) = u_2^t D_p^{-1} S_r u_2 - \beta_1 (u_2^t D_p^{-1} u_2 - 1) - \beta_2 (u_2^t D_p^{-1} u_1)$$

Entonces,

$$\frac{\partial L}{\partial u_2} = 2 D_p^{-1} S_r u_2 - 2 \beta_1 D_p^{-1} u_2 - \beta_2 D_p^{-1} u_1 = 0_{(p)}$$

$$\Rightarrow 2 u_1^t \underbrace{D_p^{-1} S_r u_2}_{0} - 2 \beta_1 \underbrace{u_1^t D_p^{-1} u_2}_{0} - \beta_2 \underbrace{u_1^t D_p^{-1} u_1}_{1} = 0_{(1)}$$

$$\text{Como } u_2^t D_p^{-1} S_r u_1 = \lambda u_2^t D_p^{-1} u_1 = 0 \Rightarrow \underline{\underline{\beta_2 = 0}}$$

Por tanto,

$$D_p^{-1} S_r u_2 = \beta_1 D_p^{-1} u_2 \Rightarrow \boxed{S_r u_2 = \beta_1 u_2}$$

y se tiene que  $(\beta_1, u_2)$  son autovalor / autovector de  $S_r$ .

• PREGUNTA 6

Demuestra que el cuadrado de la distancia euclídea entre las proyecciones de dos filas  $P_{r_i}$  y  $P_{r_{i'}}$  coincide con la distancia chi-cuadrado entre las filas  $r_i$  y  $r_{i'}$ .

$$\begin{aligned}
 d_e^2(P_{r_i}, P_{r_{i'}}) &= (P_{r_i} - P_{r_{i'}})^t (P_{r_i} - P_{r_{i'}}) = \\
 &= \sum_{\alpha=1}^P (r_i^t D_p^{-1} u_{\alpha} - r_{i'}^t D_p^{-1} u_{\alpha}) (r_i^t D_p^{-1} u_{\alpha} - r_{i'}^t D_p^{-1} u_{\alpha}) \\
 &= (r_i - r_{i'})^t D_p^{-1} \underbrace{\sum_{\alpha=1}^P u_{\alpha} u_{\alpha}^t}_{I_p} D_p^{-1} (r_i - r_{i'}) \\
 &= (r_i - r_{i'})^t D_p^{-1} (r_i - r_{i'}) = d_{\chi^2}(r_i, r_{i'})
 \end{aligned}$$

\* NOTA :

Los nuevos ejes  $u_1, \dots, u_p$  verifican

$$\begin{aligned}
 \begin{pmatrix} u_1^t \\ \vdots \\ u_p^t \end{pmatrix} D_p^{-1} (u_1 \dots u_p) &= I_p \\
 (u_1 \dots u_p) \begin{pmatrix} u_1^t \\ \vdots \\ u_p^t \end{pmatrix} D_p^{-1} &= \sum_{\alpha=1}^P u_{\alpha} u_{\alpha}^t D_p^{-1} = I_p
 \end{aligned}$$

Las coordenadas de  $r_i$  sobre los nuevos ejes son:

$$P_{r_i} : (r_i^t D_p^{-1} u_1, \dots, r_i^t D_p^{-1} u_p)$$

con

$$\begin{aligned}
 r_i &= \sum_{\alpha=1}^P (r_i^t D_p^{-1} u_{\alpha}) u_{\alpha} = \sum_{\alpha=1}^P u_{\alpha} (u_{\alpha}^t D_p^{-1} r_i) = \\
 &= \underbrace{\sum_{\alpha=1}^P u_{\alpha} u_{\alpha}^t}_{I_p} D_p^{-1} r_i
 \end{aligned}$$

• PREGUNTA 9

Demuestra que  $(0, m_r)$  es un autovalor-vector de la matriz

$$S_r = \sum_{i=1}^n f_i \cdot (r_i - m_r)(r_i - m_r)^t D_p^{-1}$$

Para ello, tenemos que comprobar que se verifica

$$S_r \cdot m_r = 0_p.$$

Luego,

$$\begin{aligned} S_r \cdot m_r &= \sum_{i=1}^n f_i \cdot (r_i - m_r)(r_i - m_r)^t \underbrace{D_p^{-1}(m_r)} = \dots \\ &= \begin{pmatrix} \frac{1}{f \cdot 1} & & \\ & \ddots & \\ & & \frac{1}{f \cdot p} \end{pmatrix} \begin{pmatrix} \frac{p \cdot 1}{f \cdot 1} \\ \vdots \\ \frac{p \cdot p}{f \cdot p} \end{pmatrix} = \\ &= \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_p = \mathbb{1}_p \end{aligned}$$

Además,

$$\left. \begin{array}{l} r_i^t \cdot \mathbb{1}_p = 1 \\ m_r^t \cdot \mathbb{1}_p = 1 \end{array} \right\} \Rightarrow (r_i - m_r)^t D_p^{-1} m_r = 0$$

Por tanto,  $S_r \cdot m_r = 0_p$

(con ello queda demostrado que  $(0, m_r)$  es un autovalor-vector de  $S_r$ ).



• PREGUNTA 10

Determina las coordenadas de las filas del autovector  $m_r$ . ¿Qué evidencian? :

$$\begin{bmatrix} r_1^t D_p^{-1} m_r \\ \vdots \\ r_n^t D_p^{-1} m_r \end{bmatrix} = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix}, \text{ donde se evidencia que la variancia de la columna es 0 (al ser la columna cte.)}$$

• PREGUNTA 11

Demuestra que la variancia ponderada de las coordenadas correspondientes al autovector

$u_k$  se puede expresar como :

$$\sum_{i=1}^n f_{i\cdot} (r_i^t D_p^{-1} u_k)^2 :$$

$$V_k = u_k^t \cdot D_p^{-1} \cdot S_r \cdot u_k = u_k^t D_p^{-1} \underbrace{\sum_{i=1}^n f_{i\cdot} (r_i - m_r)(r_i - m_r)^t D_p^{-1} u_k}_{S_r \text{ (no depende de la base)}}$$

$$= \sum_{i=1}^n f_{i\cdot} \underbrace{(r_i^t D_p^{-1} u_k - m_r^t D_p^{-1} u_k)}_0^2 = \sum_{i=1}^n f_{i\cdot} (r_i^t D_p^{-1} u_k)^2$$

" si  $u_k \neq m_r$

(Media ponderada  $m_r^t D_p^{-1} u_k = 0$  si  $u_k \neq m_r$ )

## • PREGUNTA 12

Define la inercia total y demuestra que es igual a la suma de los autovalores de  $S_r$ :

La inercia total se define como:  $InT = \frac{1}{N} X^2$

Veamos que  $InT = \frac{1}{N} X^2 = \text{tr}(S_r) = \lambda_1 + \dots + \lambda_p$

Tenemos:  $X^2 = N \cdot \sum_{i=1}^n f_i \cdot d_{X^2}(r_i, m_r) = N \cdot \sum_{i=1}^n f_i \cdot (r_i - m_r)^t D_p^{-1} (r_i - m_r)$

$$S_r = \sum_{i=1}^n f_i \cdot (r_i - m_r)(r_i - m_r)^t D_p^{-1}$$

$$\text{tr}(S_r) = \lambda_1 + \dots + \lambda_p = \sum_{i=1}^n f_i \cdot \text{tr}((r_i - m_r)(r_i - m_r)^t D_p^{-1}) =$$

$$= \sum_{i=1}^n f_i \cdot \underbrace{((r_i - m_r)^t D_p^{-1} (r_i - m_r))}_{1 \times 1 \text{ (escalares)}} = \sum_{i=1}^n f_i \cdot d_{X^2}(r_i, m_r)$$

$$\text{Entonces, } InT = \frac{1}{N} X^2 = \frac{1}{N} \cdot N \cdot \sum_{i=1}^n f_i \cdot d_{X^2}(r_i, m_r) = \text{tr}(S_r) = \lambda_1 + \dots + \lambda_p$$



## 4. ANÁLISIS DISCRIMINANTE

### • PREGUNTA 1.

Discriminación en dos poblaciones con distribuciones conocidas

(a) Define la probabilidad total de error de clasificación de una regla discriminante:

Sea  $P$  una población formada por dos grupos  $G_1$  y  $G_2$  con proporciones  $\pi_1$  y  $\pi_2$  respectivamente ( $\pi_1 + \pi_2 = 1$ )

Sea  $\underline{X} = (X_1, \dots, X_p)'$  con función de densidad  $f_i(\underline{x})$  en  $G_i$ ,

$i = 1, 2$ . Una regla discriminante viene determinada por una partición del espacio muestral  $\Omega$  en dos subconjuntos  $\Omega_1$  y  $\Omega_2$ , donde:

$$\Omega_1 \cup \Omega_2 = \Omega \quad \text{y} \quad \Omega_1 \cap \Omega_2 = \emptyset, \text{ de forma que dado}$$

un individuo  $\underline{x}_0 = (x_{01}, \dots, x_{0p})'$ , lo asignamos al grupo  $G_i$  si  $\underline{x}_0 \in \Omega_i$ ,  $i = 1, 2$ .

Las probabilidades de los errores de clasificación son:

$$P(2|1) = P(\text{Asignar en } G_2 \text{ siendo de } G_1) = \int_{\Omega_2} f_1(\underline{x}) d\underline{x}$$

$$P(1|2) = P(\text{Asignar en } G_1 \text{ siendo de } G_2) = \int_{\Omega_1} f_2(\underline{x}) d\underline{x}$$

Se define la probabilidad total del error de clasificación de una regla discriminante como:

$$P(\Omega:f) = P(\text{Asignar en } G_2 \text{ siendo de } G_1 \text{ ó asignar en } G_1 \text{ siendo de } G_2) = P(2|1) \cdot \pi_1 + P(1|2) \cdot \pi_2$$

(b) Determina la regla discriminante que minimiza la probabilidad total de error de clasificación.

$$\begin{aligned} P(\Omega : f) &= P(1|2) \cdot \pi_2 + P(2|1) \cdot \pi_1 = \pi_2 \cdot \int_{\Omega_1} f_2(x) dx + \pi_1 \int_{\Omega_2} f_1(x) dx = \\ &= \pi_2 \cdot \int_{\Omega_1} f_2(x) dx + \pi_1 \left[ 1 - \int_{\Omega_1} f_1(x) dx \right] = \\ &= \int_{\Omega_1} (\pi_2 f_2(x) - \pi_1 f_1(x)) dx + \pi_1 \end{aligned}$$

Como  $\pi_1 > 0$ ,  $P(\Omega : f)$  se minimiza donde lo haga la integral, es decir, el  $\pi_1$  final no influye en el mínimo.

Luego el recinto donde se minimiza la integral es:

$$\Omega_1^* = \left\{ x / f_2(x) \pi_2 - f_1(x) \pi_1 \leq 0 \right\} = \left\{ x / f_1(x) \pi_1 \geq f_2(x) \pi_2 \right\}$$

Así, clasificaremos un individuo en  $G_1$  si:  $\pi_1 f_1(x_0) > f_2(x_0) \pi_2$ , y en

$G_2$  si  $\pi_2 f_2(x_0) > \pi_1 f_1(x_0)$ .

(c) Plantea y resuelve el problema original de Fisher:

Sea  $P$  una población formada por dos grupos  $G_1$  y  $G_2$ . Sea

$\underline{X} = (\underline{X}_1, \dots, \underline{X}_p)'$  con distribución  $(\underline{\mu}_1, \Sigma)$  en  $G_1$  y  $(\underline{\mu}_2, \Sigma)$  en  $G_2$ .

Se plantea el problema de determinar  $b \in \mathbb{R}^p$  tal que:

$$\begin{aligned} \max_{\underline{b}} \quad & (\underline{b}'\underline{\mu}_1 - \underline{b}'\underline{\mu}_2)^2, \text{ donde } \underline{b} \text{ es el vector que más diferencia} \\ \text{s.a.: } & \underline{b}'\hat{\Sigma}\underline{b} = 1 \text{ entre los grupos.} \end{aligned}$$

Reescribimos el problema como:

$$\begin{aligned} \max_{\underline{b}} \quad & \underbrace{\underline{b}'(\underline{\mu}_1 - \underline{\mu}_2)}_d \underbrace{(\underline{\mu}_1 - \underline{\mu}_2)'}_{d'} \underline{b} \\ \text{s.a.: } & \underline{b}'\hat{\Sigma}\underline{b} = 1 \end{aligned}, \text{ donde la solución a este tipo de}$$

problema ya sabemos que viene dada por el autovector asociado al máximo autovalor de  $\hat{\Sigma}^{-1} d d^t$ , que esta matriz (en particular) solo posee un autovalor no nulo, el cual es:

$$\underline{b} = \hat{\Sigma}^{-1} \cdot d, \text{ pero no está normalizado, luego no cumple que } \underline{b}'\hat{\Sigma}\underline{b} = 1,$$

por tanto el autovector normalizado que necesitamos es, llamando

$$a = d^t \hat{\Sigma}^{-1} d, \text{ al autovalor asociado a } \underline{b}, \text{ obtenemos:}$$

$$\underline{b}^* = \frac{c}{a^{1/2}} = \frac{\hat{\Sigma}^{-1} \cdot d}{(d^t \hat{\Sigma}^{-1} d)^{1/2}}$$

(d) Demuestra que supuesto que las poblaciones se distribuyen según  $N(\mu_i, \Sigma)$ ,  $i=1,2$ , la regla del discriminante que minimiza la probabilidad total de error coincide con la propuesta por Fisher.

• Regla propuesta por Fisher:

Asignar  $x_0$  a  $G_1$  si  $|\underline{b}'x_0 - \underline{b}'\mu_1| < |\underline{b}'x_0 - \underline{b}'\mu_2| \Leftrightarrow$

$$\Leftrightarrow (\underbrace{\underline{b}'x_0}_{y_0} - \underline{b}'\mu_1)^2 < (\underbrace{\underline{b}'x_0}_{y_0} - \underline{b}'\mu_2)^2 \Leftrightarrow$$

$$\Leftrightarrow y_0^2 + (\underline{b}'\mu_1)^2 - 2y_0\underline{b}'\mu_1 < y_0^2 + (\underline{b}'\mu_2)^2 - 2y_0\underline{b}'\mu_2 \Leftrightarrow$$

$$\Leftrightarrow 2y_0 \underbrace{(\underline{b}'\mu_2 - \underline{b}'\mu_1)}_{<0} < \underbrace{(\underline{b}'\mu_2)^2 - (\underline{b}'\mu_1)^2}_{(\underline{b}'\mu_2 + \underline{b}'\mu_1)(\underline{b}'\mu_2 - \underline{b}'\mu_1)} \Leftrightarrow$$

$$\Leftrightarrow 2y_0 > \underline{b}'\mu_1 + \underline{b}'\mu_2 \Leftrightarrow \boxed{y_0 > \frac{\underline{b}'(\mu_1 + \mu_2)}{2}}$$

• Regla que minimiza la probab. del error:

Asignamos  $x_0$  a  $G_1$  si  $\pi_1 f_1(x_0) > \pi_2 f_2(x_0)$ . Para el caso particular  $f_i \sim N_p(\mu_i, \Sigma)$ ,  $i=1,2$ , si  $\pi_1 = \pi_2$ , asigno  $x_0$  en  $G_1$  si  $f_1(x_0) > f_2(x_0)$  (lo asignamos donde haya mas densidad).

Luego asigno  $x_0$  en  $G_1$  si:

Clasifico $x_0$ en $G_2$		Clasifico $x_0$ en $G_1$
$\underline{b}'\mu_2$		$\underline{b}'\mu_1$
	$\underline{b}'(\mu_1 + \mu_2)$	
	$2$	

luego asigno  $x_0$  en  $G_1$  si:

$$\boxed{\underline{b}'x_0 > \frac{\underline{b}'(\mu_1 + \mu_2)}{2}}$$

## • PREGUNTA 2

Coordenadas discriminantes canónicas.

(a) Planteamiento y solución:

Sea  $P$  una población formada por  $k$  grupos  $G_1, \dots, G_k$  con proporciones  $\pi_1, \dots, \pi_k$  respect.  $(\sum_{i=1}^k \pi_i = 1)$ . Sea  $\underline{X} = (X_1, \dots, X_p)^t$  con fdd  $f_i(x)$  en  $G_i$ ,  $i=1, \dots, k$ . Consideramos una muestra aleatoria de cada grupo  $x_{i1}, \dots, x_{ini}$  de  $G_i$ ,  $i=1, \dots, k$ .

El objetivo es determinar las comb. lineales de las vbles.  $\underline{Y} = \underline{b}^t \underline{X}$ ,  $\underline{b} \in \mathbb{R}^p$  que mas discrimina entre los grupos. Para ello consideramos el contraste:

$H_0: \underline{b}' \underline{\mu}_1 = \dots = \underline{b}' \underline{\mu}_k$ , bajo normalidad y homocedasticidad

supuestas. El estadístico es:

$$F = \alpha \frac{\underline{b}' B \underline{b}}{\underline{b}' W \underline{b}}, \text{ donde } B = \sum_{i=1}^k n_i (\underline{\bar{X}}_i - \underline{\bar{X}})(\underline{\bar{X}}_i - \underline{\bar{X}})^t, \text{ y } W = \sum_{i=1}^k S_i$$

y la región crítica  $F \geq F_\alpha$  ( $F$  crítica).

Así, debemos tomar  $\underline{b}$  de forma que el estadístico sea máximo para que las medias sean lo mas distintas posibles. Planteamos los siguientes problemas de optimización:

$$(1) \sup_{\underline{b} \neq 0} \frac{\underline{b}' B \underline{b}}{\underline{b}' W \underline{b}} = \sup_{\underline{b} \neq 0} \underline{b}' B \underline{b} \quad \text{s.a.: } \underline{b}' W \underline{b} = 1, \text{ donde la solución}$$

es el autovector  $\underline{e}_1$  asociado al máximo autovalor de  $W^{-1}B$ , y definimos la primera coordenada discriminante como  $\underline{Y}_1 = \underline{e}_1' \underline{X}$

$$(2) \sup_{\underline{b} \neq 0} \frac{\underline{b}' B \underline{b}}{\underline{b}' W \underline{b}} = \sup_{\substack{\underline{b}' B \underline{b} \\ \text{s.a.: } \underline{b}' W \underline{b} = 1}} \underline{b}' B \underline{b}, \text{ cuya solución } \underline{e}_2 \\ \underline{e}_2' W \underline{e}_2 = 0$$

es el autovector asociado al segundo mayor autovector de  $W^{-1}B$ , luego definimos la segunda coord. discriminante

$$\text{como: } Y_2 = \underline{e}_2' \underline{X}$$

$\vdots$

(Así con los  $r$  autovalores no nulos de  $W^{-1}B$ ).

(b) Describe el criterio de clasificación basado en las coord. canónicas:

Determinados  $\underline{e}_1, \dots, \underline{e}_r \in \mathbb{R}^p$ , usamos el criterio de clasificación para asignar un individuo a su grupo correspondiente.  $\underline{x}_0$  será clasificado en  $G_i$  si:

$$d(c_r \cdot \underline{x}_0, c_r \cdot \bar{\underline{x}}_i) = \min_{s=1, \dots, k} d(c_r \underline{x}_0, c_r \bar{\underline{x}}_s), \text{ donde}$$

$$c_r = (\underline{e}_1, \dots, \underline{e}_r)^t$$