

```
○○○○○
○○○○○○○○○○○○○○
○○○○○○○○○○
```

```
○○○○○
○○○○○
```

```
○○○
○○○○○○○○○
```

```
○○○○○○○
○○○○○○○○○
○○○○○
```

```
○○
○○○○○○○
○○○
```

Coordenadas Discriminantes y Discriminación

Graciela Boente



Ejemplo

Variables medidas sobre árboles de manzana de 6 injertos. Para cada injerto hay 8 árboles. Las variables son:

x_1 =Diámetro del tronco a los 4 años en unidades de 10cm,

x_2 =Largo a los 4 años,

x_3 =Diámetro del tronco a los 15 años en unidades de 10cm,

x_4 =Peso del árbol a los 15 años, en unidades de 1000 libras.

Inj.	1	1	1	1	1	1	1	1	2	2	2	2
x_1	1.11	1.19	1.09	1.25	1.11	1.08	1.11	1.16	1.05	1.17	1.11	1.25
x_2	2.569	2.928	2.865	3.844	3.027	2.336	3.211	3.037	2.074	2.885	3.378	3.906
x_3	3.58	3.75	3.93	3.94	3.60	3.51	3.98	3.62	4.09	4.06	4.87	4.98
x_4	0.760	0.821	0.928	1.009	0.766	0.726	1.209	0.750	1.036	1.094	1.635	1.517
Inj.	2	2	2	2	3	3	3	3	3	3	3	3
x_1	1.17	1.15	1.17	1.19	1.07	0.99	1.06	1.02	1.15	1.20	1.20	1.17
x_2	2.782	3.018	3.383	3.447	2.505	2.315	2.667	2.390	3.021	3.085	3.308	3.231
x_3	4.38	4.65	4.69	4.40	3.76	4.44	4.38	4.67	4.48	4.78	4.57	4.56
x_4	1.197	1.244	1.495	1.026	0.912	1.398	1.197	1.613	1.476	1.571	1.506	1.458
Inj.	4	4	4	4	4	4	4	4	5	5	5	5
x_1	1.22	1.03	1.14	1.01	0.99	1.11	1.20	1.08	0.91	1.15	1.14	1.05
x_2	2.838	2.351	3.001	2.439	2.199	3.318	3.601	3.291	1.532	2.552	3.083	2.330
x_3	3.89	4.05	4.05	3.92	3.27	3.95	4.27	3.85	4.04	4.16	4.79	4.42
x_4	0.944	1.241	1.023	1.067	0.693	1.085	1.242	1.017	1.084	1.151	1.381	1.242
Inj.	5	5	5	5	6	6	6	6	6	6	6	6
x_1	0.99	1.22	1.05	1.13	1.11	0.75	1.05	1.02	1.05	1.07	1.13	1.11
x_2	2.079	3.366	2.416	3.100	2.813	0.840	2.199	2.132	1.949	2.251	3.064	2.469
x_3	3.47	4.41	4.64	4.57	3.76	3.14	3.75	3.99	3.34	3.21	3.63	3.95
x_4	0.673	1.137	1.455	1.325	0.800	0.606	0.790	0.853	0.610	0.562	0.707	0.952



Ejemplo

Hemos estudiado si las medias de los Injertos 1, 2 y 3 eran iguales, o sea, testeamos $H_0 : \mu_1 = \mu_2 = \mu_3$.

Para ello supusimos que las matrices de covarianza eran iguales.

$$\bar{x}_1 = (1.1375, 2.9771, 3.7388, 0.8711)^T$$

$$\bar{x}_2 = (1.1575, 3.1091, 4.5150, 1.2805)^T$$

$$\bar{x}_3 = (1.1075, 2.8152, 4.4550, 1.3914)^T$$

$$\bar{x} = (1.1342, 2.9672, 4.2363, 1.1810)^T$$

Las matrices de covarianza estimadas son

$$s_1 = \begin{pmatrix} 0.0034 & 0.0203 & 0.0037 & 0.0018 \\ 0.0203 & 0.2007 & 0.0580 & 0.0458 \\ 0.0037 & 0.0580 & 0.0352 & 0.0285 \\ 0.0018 & 0.0458 & 0.0285 & 0.0283 \end{pmatrix} \quad s_2 = \begin{pmatrix} 0.0034 & 0.0258 & 0.0088 & 0.0032 \\ 0.0258 & 0.3048 & 0.1498 & 0.0832 \\ 0.0088 & 0.1498 & 0.1157 & 0.0711 \\ 0.0032 & 0.0832 & 0.0711 & 0.0565 \end{pmatrix}$$

$$s_3 = \begin{pmatrix} 0.0068 & 0.0314 & 0.0087 & 0.0060 \\ 0.0314 & 0.1543 & 0.0480 & 0.0329 \\ 0.0087 & 0.0480 & 0.0951 & 0.0680 \\ 0.0060 & 0.0329 & 0.0680 & 0.0534 \end{pmatrix}$$



Ejemplo

Construimos

$$\mathbf{U} = \mathbf{Q}_1 + \mathbf{Q}_2 + \mathbf{Q}_3 = \sum_{i=1}^3 \sum_{j=1}^{n_i} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)(\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)^T$$

$$\mathbf{H} = \sum_{i=1}^k n_i (\bar{\mathbf{x}}_i - \bar{\mathbf{x}})(\bar{\mathbf{x}}_i - \bar{\mathbf{x}})^T$$

obteniendo

$$\mathbf{U} = \begin{pmatrix} 0.0956 & 0.5422 & 0.1490 & 0.0771 \\ 0.5422 & 4.6184 & 1.7911 & 1.1326 \\ 0.1490 & 1.7911 & 1.7221 & 1.1731 \\ 0.0771 & 1.1326 & 1.1731 & 0.9670 \end{pmatrix} \quad \mathbf{H} = \begin{pmatrix} 0.0101 & 0.0592 & -0.0079 & -0.0346 \\ 0.0592 & 0.3466 & 0.0111 & -0.1674 \\ -0.0079 & 0.0111 & 2.9845 & 1.8233 \\ -0.0346 & -0.1674 & 1.8233 & 1.2014 \end{pmatrix}$$

El estadístico para testear $H_0 : \mu_1 = \mu_2 = \mu_3$ era

$$V = \frac{|\mathbf{U}|}{|\mathbf{U} + \mathbf{H}|} = \Lambda(23, 4, 2)$$

y rechazamos H_0 .



Ejemplo

En el caso de dos poblaciones, nosotros vimos que la dirección a la que se atribuía la responsabilidad del rechazo era

$$\hat{\alpha} = \mathbf{S}^{-1}(\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

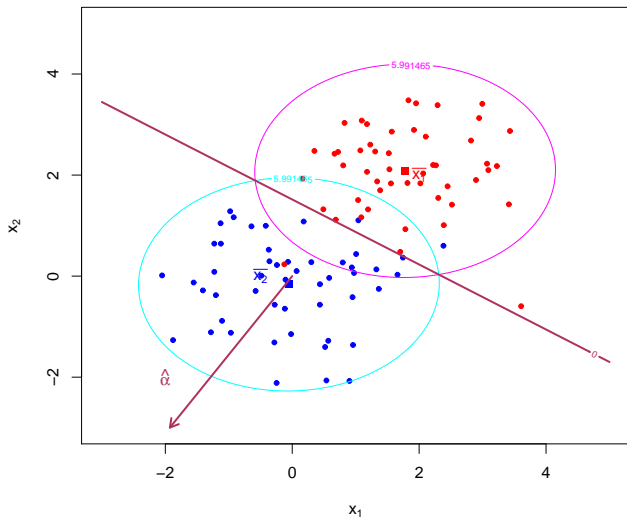
con $\mathbf{S} = ((n_1 - 1)\mathbf{S}_1 + (n_2 - 1)\mathbf{S}_2)/(n_1 + n_2 - 2)$

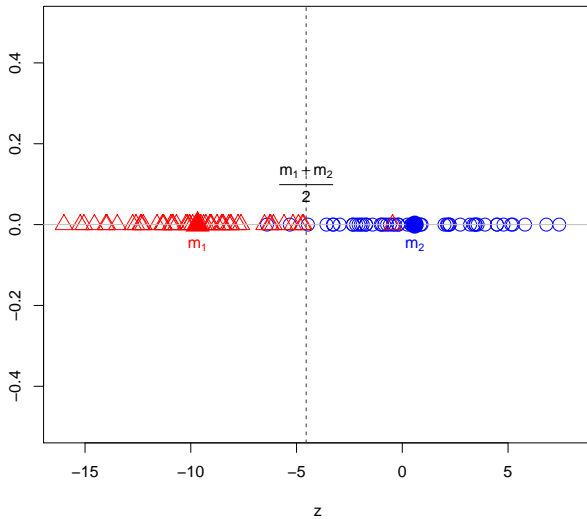
- La función $H(\mathbf{u}) = \hat{\alpha}^T \mathbf{u}$ es la función discriminante lineal.
- Asignábamos \mathbf{u} a la población 1 si

$$z = \hat{\alpha}^T \mathbf{u} > \hat{\alpha}^T \left(\frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2} \right) = \frac{m_1 + m_2}{2}$$

o equivalentemente, si $z - m_2 > m_1 - z$

Observemos que $m_1 - m_2 > 0$, luego, esto es análogo a clasificar \mathbf{u} en aquella población donde la distancia $|z - m_j|$ sea mínima.







Propiedad

Sea $\mathbf{x} \in \mathbb{R}^p$ un vector aleatorio y G una variable aleatoria que indica la pertenencia al grupo, tales que para $1 \leq i \leq k$

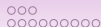
$$\mathbb{P}(G = i) = \pi_i \quad \mathbb{E}(\mathbf{x}|G = i) = \boldsymbol{\mu}_i \quad \text{VAR}(\mathbf{x}|G = i) = \boldsymbol{\Sigma}_i$$

entonces si $\bar{\boldsymbol{\mu}} = \mathbb{E}(\mathbf{x})$ y $\bar{\boldsymbol{\Sigma}} = \text{VAR}(\mathbf{x})$ se cumple que

$$\bar{\boldsymbol{\mu}} = \sum_{j=1}^k \pi_j \boldsymbol{\mu}_j \quad \bar{\boldsymbol{\Sigma}} = \boldsymbol{\Sigma}_W + \boldsymbol{\Sigma}_B$$

donde

$$\boldsymbol{\Sigma}_W = \sum_{i=1}^k \pi_i \boldsymbol{\Sigma}_i \quad \boldsymbol{\Sigma}_B = \sum_{i=1}^k \pi_i (\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})(\boldsymbol{\mu}_i - \bar{\boldsymbol{\mu}})^T$$



El Problema

Sea $z = \mathbf{a}^T \mathbf{x}$, luego

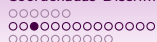
$$\text{VAR}(z) = \mathbf{a}^T \text{VAR}(\mathbf{x}) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma}_W \mathbf{a} + \mathbf{a}^T \boldsymbol{\Sigma}_B \mathbf{a}$$

es decir, descompusimos a la varianza de z en una componente que mide la variabilidad *dentro* de grupos y otra que mide la variabilidad *entre* grupos.

Nos interesan combinaciones lineales, o sea, vectores \mathbf{a} tales que la varianza de z es mucho más grande que la varianza *dentro* de grupos ya que de esto indica que la variabilidad dentro de grupos se ve aumentada por diferencias en posición.

Como $\bar{\boldsymbol{\Sigma}} = \boldsymbol{\Sigma}_W + \boldsymbol{\Sigma}_B$, nos interesarán direcciones \mathbf{a} que maximizan

$$F_{\mathbf{a}} = \frac{\mathbf{a}^T \boldsymbol{\Sigma}_B \mathbf{a}}{\mathbf{a}^T \boldsymbol{\Sigma}_W \mathbf{a}}$$



Caso $k = 2$

Si $k = 2$

$$\Sigma_B = \pi_1 \pi_2 (\mu_1 - \mu_2)(\mu_1 - \mu_2)^T$$

luego basta maximizar

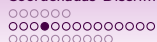
$$\frac{(\mathbf{a}^T (\mu_1 - \mu_2))^2}{\mathbf{a}^T \Sigma_W \mathbf{a}}$$

y vimos que el máximo se alcanza en

$$\alpha = \Sigma_W^{-1} (\mu_1 - \mu_2)$$

que se estima por

$$\hat{\alpha} = \mathbf{S}^{-1} (\hat{\mu}_1 - \hat{\mu}_2)$$



Caso $k > 2$

Supongamos $\Sigma_W > 0$. Sea \mathbf{C} triangular tal que $\Sigma_W = \mathbf{C}^T \mathbf{C}$ y definamos

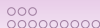
$$\mathbf{B} = (\mathbf{C}^{-1})^T \Sigma_B \mathbf{C}^{-1}$$

Sean β_1, \dots, β_p los autovectores de \mathbf{B} asociados a los autovalores $\lambda_1 \geq \dots \geq \lambda_p$.

$$\mathbf{B} = \sum_{j=1}^p \lambda_j \beta_j \beta_j^T \quad \beta_j^T \beta_\ell = 0 \text{ si } j \neq \ell \quad \|\beta_j\| = 1$$

Como $\text{rango}(\mathbf{B}) = s \leq \min(k-1, p)$, $\lambda_j = 0$ si $j > s$. Por simplicidad supondremos que los autovalores no nulos son distintos, o sea, $\lambda_1 > \dots > \lambda_s$. Es decir, tenemos que

$$\mathbf{B} = \sum_{j=1}^s \lambda_j \beta_j \beta_j^T, \quad \beta_j^T \beta_\ell = 0 \text{ si } 1 \leq j \neq \ell \leq p \quad \|\beta_j\| = 1, 1 \leq j \leq p$$



Caso $k > 2$

Luego,

$$F_{\mathbf{a}} = \frac{\mathbf{a}^T \boldsymbol{\Sigma}_B \mathbf{a}}{\mathbf{a}^T \boldsymbol{\Sigma}_W \mathbf{a}} = \frac{(\mathbf{C}\mathbf{a})^T \mathbf{B} (\mathbf{C}\mathbf{a})}{\|\mathbf{C}\mathbf{a}\|^2}$$

Por lo tanto, si $\mathbf{b} = \mathbf{C}\mathbf{a}$

$$\max_{\mathbf{a} \neq \mathbf{0}} F_{\mathbf{a}} = \max_{\mathbf{b} \neq \mathbf{0}} \frac{\mathbf{b}^T \mathbf{B} \mathbf{b}}{\|\mathbf{b}\|^2}$$

El máximo de la expresión del lado derecho se alcanza en β_1 el autovector de \mathbf{B} asociado a su mayor autovalor λ_1 .

Con lo cual,

$$\max_{\mathbf{a} \neq \mathbf{0}} F_{\mathbf{a}} = F_{\alpha_1} = \lambda_1$$

donde

$$\alpha_1 = \mathbf{C}^{-1} \beta_1$$

La combinación lineal $z_1 = \alpha_1^T \mathbf{x}$ se llama la primer **coordenada discriminante** y da la mejor separación entre grupos.



Caso $k > 2$

Tenemos que si rechazamos H_3

$$0 < \text{rango}(\mathbf{B}) = \text{rango}(\mathbf{\Sigma}_B) = s \leq \min(k - 1, p).$$

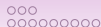
Luego, para elegir las siguientes direcciones y no repetir información, buscaremos maximizar F_a sujeto a la condición de que las nuevas coordenadas sean no-correlacionadas con z_1 , es decir, tales que

$$\text{Cov}(\mathbf{a}^T \mathbf{x}, \alpha_1^T \mathbf{x}) = \mathbf{a}^T \mathbf{\Sigma}_W \alpha_1 = 0.$$

Este problema se puede escribir como

$$\max_{\substack{\mathbf{a} \neq 0 \\ \mathbf{a}^T \mathbf{\Sigma}_W \alpha_1 = 0}} F_a = \max_{\substack{\mathbf{b} \neq 0 \\ \mathbf{b}^T \beta_1 = 0}} \frac{\mathbf{b}^T \mathbf{B} \mathbf{b}}{\|\mathbf{b}\|^2}$$

El lado derecho de la expresión se alcanza en β_2 el autovector de \mathbf{B} asociado al segundo mayor autovalor λ_2 .



Caso $k > 2$

$$\max_{\mathbf{a} \neq \mathbf{0}} F_{\mathbf{a}} = F_{\alpha_2} = \lambda_2 \quad \alpha_2 = \mathbf{C}^{-1}\beta_2$$

$$\mathbf{a}^T \boldsymbol{\Sigma}_W \alpha_1 = 0$$

donde β_2 el autovector de \mathbf{B} asociado a λ_2 .

Definición:

Definimos la j -ésima variable canónica o variable discriminante z_j como

$$z_j = \alpha_j^T \mathbf{x} \quad \text{donde} \quad \alpha_j = \mathbf{C}^{-1}\beta_j$$

El vector $\mathbf{z} = (z_1, \dots, z_p)^T$ es el vector de variables canónicas o variables discriminantes.

$$\mathbf{z} = \mathbf{A}^T \mathbf{x} \quad \text{con} \quad \mathbf{A} = (\alpha_1, \dots, \alpha_p)$$

Observemos que α_j es un autovector de $\boldsymbol{\Sigma}_W^{-1} \boldsymbol{\Sigma}_B$.



Caso $k > 2$ y $s < p$

Como $\text{rango}(\mathbf{B}) = s \leq \min(k - 1, p)$ si $s < p$, para todo $s + 1 \leq m \leq p$ tendremos que

$$\max_{\mathbf{a} \neq \mathbf{0}} F_{\mathbf{a}} = F_{\alpha_m} = \lambda_m = 0$$

$$\mathbf{a}^T \boldsymbol{\Sigma}_W \boldsymbol{\alpha}_\ell = 0, \ell < m$$

Es decir, una combinación lineal no correlacionada con z_1, \dots, z_s no dará información sobre las diferencias en posición.

En particular, z_{s+1}, \dots, z_p no dan información sobre las diferencias de posición.

O sea, si $\bar{\nu} = \mathbf{A}^T \bar{\mu}$ y si $\nu_i = \mathbf{A}^T \mu_i$, entonces

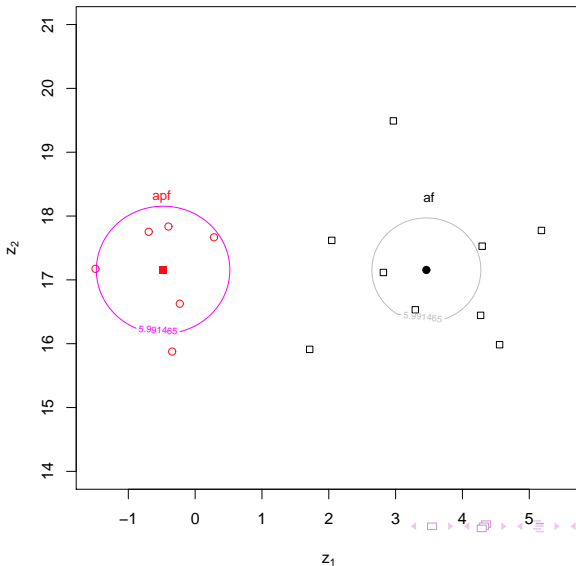
$$\nu_{im} = \bar{\nu}_m \quad \text{para} \quad s + 1 \leq m \leq p$$

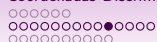
Más aún, si $\boldsymbol{\Lambda}_1 = \text{diag}(\lambda_1, \dots, \lambda_s)$

$$\mathbf{A}^T \boldsymbol{\Sigma}_B \mathbf{A} = \begin{pmatrix} \boldsymbol{\Lambda}_1 & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{pmatrix}$$



Ejemplo





Caso $k > 2$ y $s < p$

Por lo tanto, sólo basta considerar el vector de variables discriminantes

$$\mathbf{z}^{(1)} = (z_1, \dots, z_s) = \mathbf{A}_1^T \mathbf{x}$$

donde

$$\mathbf{A} = (\mathbf{A}_1, \mathbf{A}_2) \quad \text{con} \quad \mathbf{A}_1 = (\boldsymbol{\alpha}_1, \dots, \boldsymbol{\alpha}_s)$$

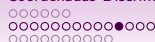
Sean además, $\boldsymbol{\nu}_i^{(1)} = \mathbf{A}_1^T \boldsymbol{\mu}_i$, $\mathbf{z}^{(2)} = \mathbf{A}_2^T \mathbf{x}$, $\mathbf{z} = \begin{pmatrix} \mathbf{z}^{(1)} \\ \mathbf{z}^{(2)} \end{pmatrix}$

Si $\boldsymbol{\Sigma}_i = \boldsymbol{\Sigma}$, $1 \leq i \leq k$ entonces $\boldsymbol{\Sigma}_W = \boldsymbol{\Sigma}$ por lo tanto,

$$\text{VAR}(\mathbf{z}|G = i) = \mathbf{I}_p \quad \mathbb{E}(\mathbf{z}|G = i) = \boldsymbol{\nu}_i$$

Si $\mathbf{x}|G = i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ entonces

$$\mathbf{z}|G = i \sim N(\boldsymbol{\nu}_i, \mathbf{I}_p)$$



Caso $k > 2$

En particular

$$\mathbf{z}^{(1)} | G = i \sim N(\boldsymbol{\nu}_i^{(1)}, \mathbf{I}_s)$$

Luego, una manera heurística de definir una región para clasificar una nueva observación es asignar el vector \mathbf{x}_0 al grupo i si la media $\boldsymbol{\nu}_i^{(1)}$ de las variables transformadas es la más cercana a $\mathbf{v}_0 = \mathbf{A}_1^T \mathbf{x}_0$.

O sea, asigno \mathbf{x}_0 al grupo i si $\mathbf{v}_0 = \mathbf{A}_1^T \mathbf{x}_0 \in \mathcal{G}_i$ donde

$$\begin{aligned} \mathcal{G}_i &= \{ \mathbf{v} \in \mathbb{R}^s : \|\mathbf{v} - \boldsymbol{\nu}_i^{(1)}\| < \|\mathbf{v} - \boldsymbol{\nu}_\ell^{(1)}\| \quad \forall \ell \neq i \} \\ &= \{ \mathbf{v} \in \mathbb{R}^s : (\boldsymbol{\nu}_i^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_i^{(1)}) > (\boldsymbol{\nu}_\ell^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_\ell^{(1)}) \quad \forall \ell \neq j \} \end{aligned}$$

Si $s = 1$ partimos la recta de las observaciones transformadas en dos semi-rectas.

Si $s = 2$ la regla de clasificación genera una partición de \mathbb{R}^2 en regiones limitadas por semi-rectas.



Inferencia

Supongamos tener ahora observaciones $\mathbf{x}_{ij} \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$,
 $1 \leq j \leq n_i, 1 \leq i \leq k$, entonces

- $\bar{\mathbf{x}}_i$ es el EMV de $\boldsymbol{\mu}_i$.
- El EMV de $\boldsymbol{\Sigma}_W$ es

$$\frac{\mathbf{U}}{n}$$

- El EMV de $\boldsymbol{\Sigma}_B$ es

$$\frac{\mathbf{H}}{n}$$

y $\mathbb{P}(\text{rango}(\mathbf{H}) = \min(p, k - 1)) = 1$

- Un estimador insesgado de $\boldsymbol{\Sigma}_W$ es

$$\frac{\mathbf{U}}{n - k}$$



Inferencia

- $\mathbf{U} = \mathbf{T}^T \mathbf{T}$
- \mathbf{b}_j los autovectores de $\hat{\mathbf{B}} = (\mathbf{T}^{-1})^T \mathbf{H} \mathbf{T}^{-1}$ tales que $\|\mathbf{b}_j\| = 1$, $\mathbf{b}_j^T \mathbf{b}_\ell = 0$ si $\ell \neq j$
 \mathbf{b}_j es el autovector asociado al j -ésimo autovalor $\hat{\lambda}_j$ de $\hat{\mathbf{B}}$, donde

$$\mathbb{P}(\hat{\lambda}_1 > \hat{\lambda}_2 > \dots > \hat{\lambda}_s > 0 \quad \text{y} \quad s = \min(p, k - 1)) = 1$$

- $\mathbf{a}_j = (\mathbf{T}^{-1})^T \mathbf{b}_j \sqrt{n - k}$, $\hat{\mathbf{A}}_1 = (\mathbf{a}_1, \dots, \mathbf{a}_s)$
- $\bar{\mathbf{z}}_i^{(1)} = \hat{\mathbf{A}}_1^T \bar{\mathbf{x}}_i$. Entonces, si n_i es grande $\bar{\mathbf{z}}_i^{(1)} \approx N(\boldsymbol{\nu}_i^{(1)}, \mathbf{I}_s/n_i)$.



Inferencia

Una región de confianza asintótica de nivel $1 - \alpha$ para $\nu_i^{(1)}$ está dada por

$$\{\mathbf{v} \in \mathbb{R}^s : n_i \|\bar{\mathbf{z}}_i^{(1)} - \mathbf{v}\|^2 \leq \chi_{s,\alpha}^2\} = \{\mathbf{v} \in \mathbb{R}^s : \|\bar{\mathbf{z}}_i^{(1)} - \mathbf{v}\| \leq \sqrt{\frac{\chi_{s,\alpha}^2}{n_i}}\}$$

Además, asignamos \mathbf{x}_0 al grupo i si $\hat{\mathbf{v}}_0 = \hat{\mathbf{A}}_1^T \mathbf{x}_0 \in \hat{\mathcal{G}}_i$ donde

$$\begin{aligned} \hat{\mathcal{G}}_i &= \{\mathbf{v} \in \mathbb{R}^s : \|\mathbf{v} - \bar{\mathbf{z}}_i^{(1)}\| < \|\mathbf{v} - \bar{\mathbf{z}}_\ell^{(1)}\| \quad \forall \ell \neq i\} \\ &= \{\mathbf{v} \in \mathbb{R}^s : (\bar{\mathbf{z}}_i^{(1)})^T (\mathbf{v} - \frac{1}{2} \bar{\mathbf{z}}_i^{(1)}) > (\bar{\mathbf{z}}_\ell^{(1)})^T (\mathbf{v} - \frac{1}{2} \bar{\mathbf{z}}_\ell^{(1)}) \quad \forall \ell \neq j\} \\ &= \{\mathbf{v} \in \mathbb{R}^s : (\bar{\mathbf{z}}_i^{(1)} - \bar{\mathbf{z}}_\ell^{(1)})^T \left(\mathbf{v} - \frac{1}{2} (\bar{\mathbf{z}}_i^{(1)} + \bar{\mathbf{z}}_\ell^{(1)}) \right) > 0 \quad \forall \ell \neq j\} \end{aligned}$$



Ejemplo

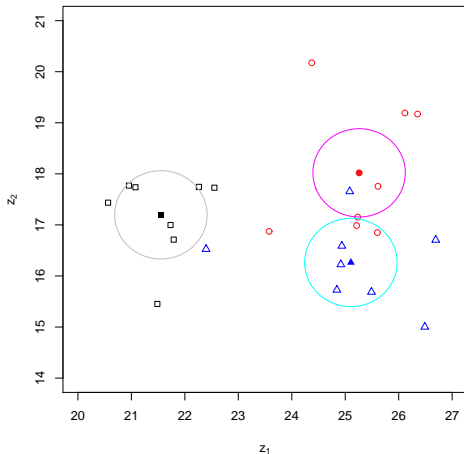
En el ejemplo de los árboles tenemos $k = 3$, $p = 4$ luego $s = 2$, los autovalores no nulos de $\hat{\mathbf{B}} = (\mathbf{T}^{-1})^T \mathbf{H} \mathbf{T}^{-1}$ son 3.3522, 0.5879 y

$$\hat{\mathbf{B}} = \begin{pmatrix} 0.0139 & -0.4232 & 0.0000 & 0.9059 \\ -0.0139 & -0.0207 & 0.9996 & -0.0094 \\ -0.9845 & 0.1523 & -0.0098 & 0.0863 \\ 0.1742 & 0.8929 & 0.0248 & 0.4145 \end{pmatrix}$$

$$\hat{\mathbf{A}}_1 = \begin{pmatrix} 7.7167 & 0.0711 \\ -2.7909 & 0.8135 \\ 6.1042 & 6.3137 \\ -1.9958 & -10.2292 \end{pmatrix}$$

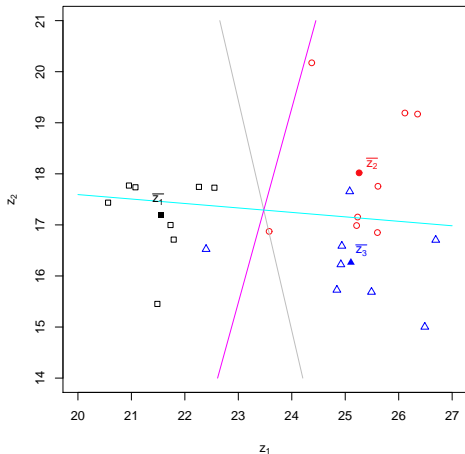


Ejemplo: Se grafican los círculos de radio $\{\chi_{2,0.05}^2/n_i\}^{1/2} = 2.4477/\sqrt{n_i}$



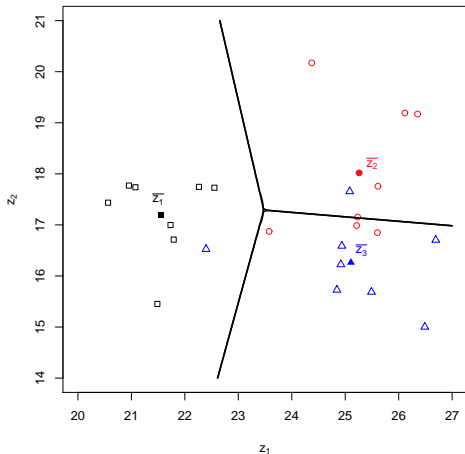


Ejemplo





Ejemplo





Ejemplo 2

Medidas tridimensionales en cráneos de 4 subespecies de oso hormiguero. Las mediciones hechas fueron

- y_1 = Largo de base, excluyendo premaxilar
- y_2 = Largo oxipitonasal
- y_3 = Largo de los nasales.

Los datos consisten en $x_j = \log(y_j)$ para las subespecies

- *Instabilis*, (Colombia) $n_1 = 21$, $\bar{x}_1 = (2.054, 2.066, 1.621)^T$.
- *Chapadensis*, en tres localidades
 - Minas Gerais: $n_2 = 6$, $\bar{x}_2 = (2.097, 2.1, 1.625)^T$,
 - Matto Grosso: $n_3 = 9$, $\bar{x}_3 = (2.091, 2.095, 1.624)^T$,
 - Santa Cruz: $n_4 = 3$, $\bar{x}_4 = (2.099, 2.102, 1.643)^T$
- *Chiriquensis*, (Panamá) $n_5 = 4$, $\bar{x}_5 = (2.092, 2.11, 1.703)^T$
- *Mexicana* $n_6 = 5$, $\bar{x}_5 = (2.099, 2.107, 1.671)^T$

Coordenadas Discriminantes



Discriminación



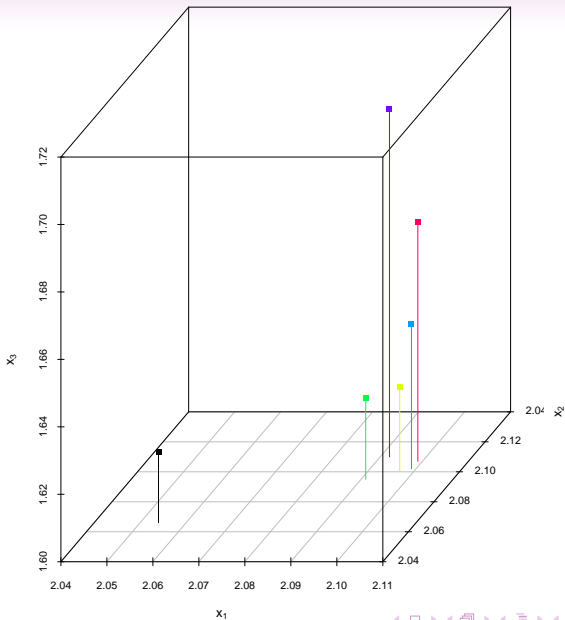
Errores de clasificación



Caso Normal I



Caso Normal II





Ejemplo 2

La matrices \mathbf{U} y \mathbf{H} son

$$\mathbf{U} = \begin{pmatrix} 0.01363 & 0.01277 & 0.01644 \\ & 0.01293 & 0.01714 \\ & & 0.03615 \end{pmatrix} \quad \mathbf{H} = \begin{pmatrix} 0.02002 & 0.01744 & 0.01308 \\ & 0.01585 & 0.01507 \\ & & 0.03068 \end{pmatrix}$$

Luego, $\nu_1 = pk + \frac{p(p+1)}{2} = 24$ y $\nu_2 = p + \frac{p(p+1)}{2} = 9$ de donde

$$-48 \log \left(\frac{|\mathbf{U}|}{|\mathbf{U} + \mathbf{H}|} \right) \approx \chi_{15}^2$$

$$\frac{|\mathbf{U}|}{|\mathbf{U} + \mathbf{H}|} = 0.1468 \quad -48 \log \left(\frac{|\mathbf{U}|}{|\mathbf{U} + \mathbf{H}|} \right) = 92.08 \quad \chi_{15,0.01}^2 = 30.57791$$

y rechazo la igualdad de medias.



Ejemplo 2

En este ejemplo $s = 3$. Los autovalores de \hat{B} son

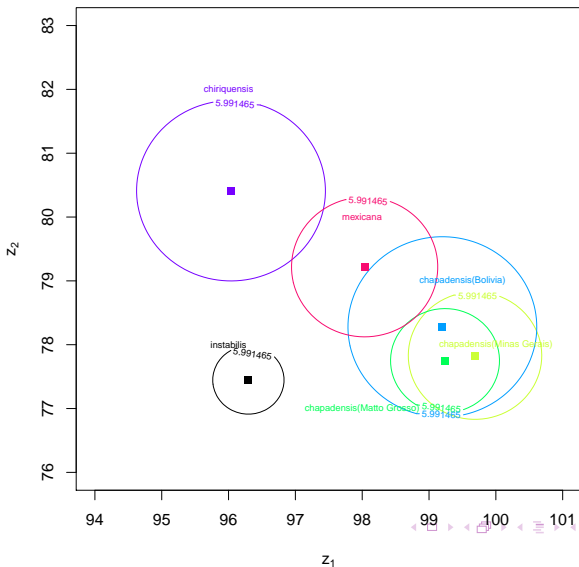
$$\hat{\lambda}_1 = 2.4001 \quad \hat{\lambda}_2 = 0.9050 \quad \hat{\lambda}_3 = 0.0515$$

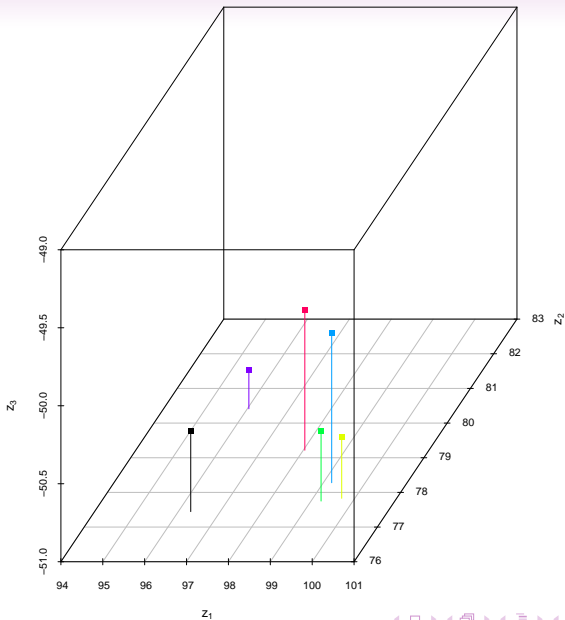
$$\hat{\mathbf{A}}_1 = \begin{pmatrix} 108.9220 & -40.6120 & 169.3528 \\ -33.8285 & 59.9577 & -222.1998 \\ -35.4962 & 22.8195 & 37.4651 \end{pmatrix}$$

Graficaremos las dos primeras coordenadas discriminantes.



Ejemplo 2: Se grafican los círculos de radio $\{\chi_{2,0.05}^2/n_i\}^{1/2} = 2.4477/\sqrt{n_i}$







Tenemos k poblaciones o grupos diferentes $\mathcal{P}_1, \dots, \mathcal{P}_k$ y un vector $\mathbf{x} \in \mathbb{R}^p$ que puede pertenecer a cualquiera de esas poblaciones.

a) G una variable aleatoria que indica la pertenencia al grupo,

$$\mathbb{P}(G = j) = \pi_j$$

Enfoque bayesiano,

- $\Theta = \{\mathcal{P}_1, \dots, \mathcal{P}_k\}$ o en forma simplificada tenemos una variable aleatoria G que toma valores $\{1, \dots, k\}$.
- Sobre Θ definimos una probabilidad a priori τ que es discreta y es tal que $\mathbb{P}(G = j) = \pi_j$

b) La distribución de \mathbf{x} varía según el grupo de pertenencia.

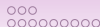
Si \mathbf{x} pertenece a la población \mathcal{P}_j entonces \mathbf{x} tiene densidad f_j ,

$$\mathbf{x} | G = j \sim f_j$$

c) Para dar una regla de clasificación daremos una partición de \mathbb{R}^p en k conjuntos disjuntos

$$\mathbb{R}^p = \bigcup_{i=1}^k \mathcal{G}_i$$

$$\mathcal{G}_j \cap \mathcal{G}_s = \emptyset$$



Definiciones

La densidad marginal de \mathbf{x} está dada por

$$f_{\mathbf{x}}(\mathbf{x}) = \sum_{j=1}^k \pi_j f_j(\mathbf{x})$$

y la probabilidad condicional de que una observación pertenezca a \mathcal{P}_j dado que $\mathbf{x} = \mathbf{x}_0$, está dada por

$$q_j(\mathbf{x}_0) = \mathbb{P}(G = j | \mathbf{x} = \mathbf{x}_0) = \frac{\pi_j f_j(\mathbf{x}_0)}{\sum_{\ell=1}^k \pi_{\ell} f_{\ell}(\mathbf{x}_0)}$$

La cantidad $q_j(\mathbf{x}_0)$ es la probabilidad a posteriori.

Definición 1.

Una regla de clasificación es una variable aleatoria $G^*(\mathbf{x})$ tal que

$$G^*(\mathbf{x}) = j \quad \text{si} \quad \mathbf{x} \in \mathcal{G}_j$$

donde $\{\mathcal{G}_1, \dots, \mathcal{G}_k\}$ es una partición de \mathbb{R}^p .



Definiciones

Podemos ver a G^* como la pertenencia predicha mientras que G es la pertenencia real.

La teoría de clasificación trata de encontrar reglas de clasificación óptimas en algún sentido. Lo ideal sería que $\mathbb{P}(G^* = G) = 1$, pero esto no es posible.

Definición 2. Para una regla de clasificación G^* con regiones de clasificación $\{\mathcal{G}_1, \dots, \mathcal{G}_k\}$, la probabilidad de asignar la observación \mathbf{x} a \mathcal{P}_i cuando en realidad, $\mathbf{x} \in \mathcal{P}_j$ es

$$p_{i|j} = \mathbb{P}(G^* = i | G = j) = \mathbb{P}(\mathbf{x} \in \mathcal{G}_i | G = j) = \int_{\mathcal{G}_i} f_j(\mathbf{x}) d\mathbf{x}$$

Observemos que $\sum_{i=1}^k p_{i|j} = 1$.



Definiciones

A veces es posible asignar un costo $c_{i|j} \geq 0$ a la clasificación de una observación del grupo j en el grupo i . En muchos casos se elige $c_{i|j} = 1$ si $i \neq j$.

Definimos

a) La función de pérdida como

$$L(\mathcal{P}_j, i) = \begin{cases} c_{i|j} & \text{si } i \neq j \\ 0 & \text{si } i = j \end{cases}$$

$$L(\mathcal{P}_j, G^*) = L(j, G^*) = \sum_{i=1}^k c_{i|j} \mathbb{I}(\mathbf{x} \in \mathcal{G}_i | \mathbf{x} \in \mathcal{P}_j)$$

donde $c_{j|j} = 0$.

b) El riesgo de G^* es

$$R(\mathcal{P}_j, G^*) = \mathbb{E} L(\mathcal{P}_j, G^*) = \sum_{i=1}^k c_{i|j} p_{i|j} = \sum_{i=1}^k c_{i|j} \int_{\mathcal{G}_i} f_j(\mathbf{x}) d\mathbf{x}$$



Definiciones

El Riesgo de Bayes de una regla de clasificación G^* será

$$r(\tau, G^*) = \mathbb{E}R(\Theta, G^*) = \sum_{j=1}^k \sum_{i \neq j} \pi_j c_{ij} p_{i|j} = \sum_{j=1}^k \sum_{i \neq j} \pi_j c_{ij} \int_{G_i} f_j(\mathbf{x}) d\mathbf{x}$$

En particular, si $c_{ij} = 1$ si $i \neq j$ tenemos que

$$r(\tau, G^*) = \sum_{j=1}^k \sum_{i \neq j} \pi_j p_{i|j} = 1 - \sum_{j=1}^k \pi_j p_{j|j} = 1 - \sum_{j=1}^k \pi_j \int_{G_j} f_j(\mathbf{x}) d\mathbf{x}$$

que se llama la **probabilidad total de mala clasificación** ya que coincide con $\mathbb{P}(G^* \neq G)$.

Si $k = 2$ y $c_{ij} = 1$ si $i \neq j$ tenemos que

$$r(\tau, G^*) = \pi_1 + \int_{G_1} [\pi_2 f_2(\mathbf{x}) - \pi_1 f_1(\mathbf{x})] d\mathbf{x}$$



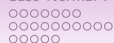
Definiciones

1. Diremos que una regla de clasificación G_0^* es Bayes respecto de la distribución a priori τ si

$$r(\tau, G_0^*) = \min_{G^*} r(\tau, G^*)$$

2. Diremos que una regla de clasificación G_0^* es minimax si

$$\max_{1 \leq j \leq k} R(\mathcal{P}_j, G_0^*) = \min_{G^*} \max_{1 \leq j \leq k} R(\mathcal{P}_j, G^*)$$



Propiedad

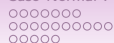
La regla Bayes respecto de τ clasifica $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \mathcal{G}_{i,0}$ donde

$$\mathcal{G}_{i,0} = \{\mathbf{x} \in \mathbb{R}^p : \sum_{j=1}^k \pi_j c_{i|j} f_j(\mathbf{x}) < \sum_{j=1}^k \pi_j c_{\ell|j} f_j(\mathbf{x}) \quad \forall \ell \neq i\}$$

o sea, clasifico $\mathbf{x} \in \mathcal{P}_i$ si

$$\sum_{j=1}^k \pi_j c_{i|j} f_j(\mathbf{x}) = \min_{\ell} \sum_{j=1}^k \pi_j c_{\ell|j} f_j(\mathbf{x})$$

siendo la asignación en la frontera de $\mathcal{G}_{i,0}$ arbitraria.



Casos particulares

- a) Supongamos que $c_{ij} = 1$ si $i \neq j$ entonces la regla Bayes clasifica $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \mathcal{G}_{i,0}$ donde

$$\begin{aligned}\mathcal{G}_{i,0} &= \{\mathbf{x} \in \mathbb{R}^P : \pi_\ell f_\ell(\mathbf{x}) < \pi_i f_i(\mathbf{x}) \quad \forall \ell \neq i\} \\ &= \{\mathbf{x} \in \mathbb{R}^P : q_\ell(\mathbf{x}) < q_i(\mathbf{x}) \quad \forall \ell \neq i\}\end{aligned}$$

es decir, clasifico $\mathbf{x} \in \mathcal{P}_i$ si $q_i(\mathbf{x}) = \max_{1 \leq \ell \leq k} q_\ell(\mathbf{x})$ siendo la asignación en la frontera de $\mathcal{G}_{i,0}$ arbitraria. Por lo tanto,

- i) la regla Bayes coincide con el criterio de minimizar la probabilidad total de mala clasificación.
- ii) la regla Bayes coincide con el criterio de maximizar la probabilidad a posteriori.
- iii) si además, $\pi_j = 1/k$, $1 \leq j \leq k$, la regla Bayes coincide con el criterio de máxima verosimilitud, que asigna \mathbf{x} a la población que maximiza la verosimilitud de \mathbf{x} .



Casos particulares

- b) Supongamos que $k = 2$ entonces la regla Bayes clasifica $\mathbf{x} \in \mathcal{P}_1$ si $\mathbf{x} \in \mathcal{G}_{1,0}$ donde

$$\mathcal{G}_{1,0} = \left\{ \mathbf{x} \in \mathbb{R}^p : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > \frac{\pi_2 c_{1|2}}{\pi_1 c_{2|1}} \right\}$$

- i) Si $\pi_2 c_{1|2} = \pi_1 c_{2|1}$ la regla Bayes da el criterio de máxima verosimilitud. En particular, si $c_{i|j} = 1$ si $i \neq j$, el criterio de máxima verosimilitud es la regla Bayes asociada a $\pi_1 = \pi_2 = \frac{1}{2}$
- ii) si $c_{i|j} = 1$ si $i \neq j$ y $\pi_1 = 1 - \alpha$, $\pi_2 = \alpha$, $0 < \alpha < 1$, entonces la regla Bayes clasifica
- $\mathbf{x} \in \mathcal{P}_1$ si

$$\mathbf{x} \in \mathcal{G}_{1,0} = \left\{ \mathbf{x} \in \mathbb{R}^p : \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} > a = \frac{\alpha}{1 - \alpha} \right\}$$

- $\mathbf{x} \in \mathcal{P}_2$ si

$$\mathbf{x} \in \mathcal{G}_{2,0} = \left\{ \mathbf{x} \in \mathbb{R}^p : \frac{f_2(\mathbf{x})}{f_1(\mathbf{x})} > \frac{1}{a} = \frac{1 - \alpha}{\alpha} \right\}$$



Casos particulares

- b) ii) Sea a_0 tal que

$$\int_{\mathcal{G}_{1,0}} f_2(\mathbf{x}) d\mathbf{x} = \int_{\mathcal{G}_{2,0}} f_1(\mathbf{x}) d\mathbf{x}$$

Luego, la regla Bayes respecto de $\tau = (\alpha_0, 1 - \alpha_0)$ con $\alpha_0 = 1/(1 + a_0)$ iguala riesgos y es la regla minimax.

- c) Supongamos que $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)^T$ con $\mathbf{x}_1 \in \mathbb{R}^q$ y que \mathbf{x}_1 y \mathbf{x}_2 son independientes en todas las poblaciones, o sea, $f_j(\mathbf{x}) = h_j(\mathbf{x}_1)\ell_j(\mathbf{x}_2)$. Más aún, supongamos que $\ell_j(\mathbf{x}_2) = \ell(\mathbf{x}_2)$ para todo j . Entonces, la regla de clasificación se basa solamente en \mathbf{x}_1 , es decir, la regla Bayes clasifica $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \mathcal{G}_{i,0}$ donde

$$\mathcal{G}_{i,0} = \{\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^p : \sum_{j=1}^k \pi_j c_{i|j} h_j(\mathbf{x}_1) < \sum_{j=1}^k \pi_j c_{\ell|j} h_j(\mathbf{x}_1) \quad \forall \ell \neq i\}$$

o sea, clasifico $\mathbf{x} \in \mathcal{P}_i$ si $\sum_{j=1}^k \pi_j c_{i|j} h_j(\mathbf{x}_1) = \min_{\ell} \sum_{j=1}^k \pi_j c_{\ell|j} h_j(\mathbf{x}_1)$.



Ejemplo

Supongamos que $k = 2$, $p = 2$

- $\mathbf{x} = (x_1, x_2)^T$ con x_i Bernouilli independientes
- $\mathbb{P}(x_1 = 1 | G = 1) = p_{11}$, $\mathbb{P}(x_1 = 1 | G = 2) = p_{21}$,
- $\mathbb{P}(x_2 = 1 | G = 1) = p_{12}$, $\mathbb{P}(x_2 = 1 | G = 2) = p_{22}$

Luego, la regla de clasificación óptima (Bayes) con costos iguales lleva a una función discriminante lineal, o sea, decido que $\mathbf{x} \in \mathcal{P}_1$ si

$$\beta_0 + \beta_1 x_1 + \beta_2 x_2 > \log \left(\frac{\pi_2}{\pi_1} \right)$$

donde

$$\beta_0 = \log \left(\frac{(1 - p_{11})(1 - p_{12})}{(1 - p_{21})(1 - p_{22})} \right)$$

$$\beta_1 = \log \left(\frac{p_{11}(1 - p_{21})}{(1 - p_{11})p_{21}} \right) \quad \beta_2 = \log \left(\frac{p_{12}(1 - p_{22})}{(1 - p_{12})p_{22}} \right)$$



Problema

Hasta ahora supusimos que la distribución de \mathbf{x} en cada población es conocida. En la mayoría de los casos esto no ocurre y tenemos alguna de las siguientes situaciones

- a) la distribución es conocida salvo por algunos parámetros que deberemos estimar, $f_j = f_j(\cdot, \theta_j)$
- b) la distribución es parcialmente desconocida, o sea, sabemos por ejemplo que

$$\log \frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} = \alpha + \beta^T \mathbf{x}$$

- c) la distribución es desconocida

En a) y b) estimamos los parámetros. La regla en este caso se estima por \hat{G}_0^* reemplazando los parámetros desconocidos por sus estimadores.

Entonces, necesitamos conocer las probabilidades de error cometido, o sea, aproximar $r(\tau, G_0^*)$ y $R(\mathcal{P}_j, G_0^*)$.



Error óptimo de clasificación

El error de mala clasificación de la población j

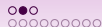
$$R(\mathcal{P}_j, G_0^*) = \sum_{i \neq j} p_{i|j} = \sum_{i \neq j} \int_{G_i} f_j(\mathbf{x}) d\mathbf{x}$$

Si $k = 2$, llamaremos

$$e_{1,opt} = R(\mathcal{P}_1, G_0^*) = \int_{G_{2,0}} f_1(\mathbf{x}) d\mathbf{x} = 1 - \int_{G_{1,0}} f_1(\mathbf{x}) d\mathbf{x}$$

$$e_{2,opt} = R(\mathcal{P}_2, G_0^*) = \int_{G_{1,0}} f_2(\mathbf{x}) d\mathbf{x} = 1 - \int_{G_{2,0}} f_2(\mathbf{x}) d\mathbf{x}$$

$$e_{opt} = \pi_1 e_{1,opt} + \pi_2 e_{2,opt}$$



Error óptimo de clasificación si $f_j = f_j(\cdot, \theta_j)$

El error de mala clasificación de la población j

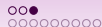
$$R(\mathcal{P}_j, G_0^*) = \sum_{i \neq j} p_{i|j} = \sum_{i \neq j} \int_{\mathcal{G}_i} f_j(\mathbf{x}, \theta_j) d\mathbf{x}$$

Si $k = 2$, llamaremos

$$e_{1,opt} = R(\mathcal{P}_1, G_0^*) = \int_{\mathcal{G}_{2,0}} f_1(\mathbf{x}, \theta_1) d\mathbf{x} = 1 - \int_{\mathcal{G}_{1,0}} f_1(\mathbf{x}, \theta_1) d\mathbf{x}$$

$$e_{2,opt} = R(\mathcal{P}_2, G_0^*) = \int_{\mathcal{G}_{1,0}} f_2(\mathbf{x}, \theta_2) d\mathbf{x} = 1 - \int_{\mathcal{G}_{2,0}} f_2(\mathbf{x}, \theta_2) d\mathbf{x}$$

$$e_{opt} = \pi_1 e_{1,opt} + \pi_2 e_{2,opt}$$



$f_j = f_j(\cdot, \theta_j)$, θ_j desconocido

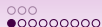
Hasta ahora supusimos que la distribución de \mathbf{x} en cada población es conocida. Supongamos que la distribución es conocida salvo por algunos parámetros que deberemos estimar, $f_j = f_j(\cdot, \theta_j)$ y sea

- $\hat{\theta}_j$ un estimador de θ_j basado en la muestra $\mathbf{x}_{j,1}, \dots, \mathbf{x}_{j,n_j}$.
- $\hat{f}_j(\mathbf{x}) = f_j(\cdot, \hat{\theta}_j)$

entonces la regla Bayes, con $c_{ij} = 1$ si $i \neq j$, se estima por la regla \hat{G}_0^* que clasifica $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \hat{G}_{i,0}$ donde

$$\hat{G}_{i,0} = \{\mathbf{x} \in \mathbb{R}^p : \pi_\ell \hat{f}_\ell(\mathbf{x}) < \pi_i \hat{f}_i(\mathbf{x}) \quad \forall \ell \neq i\}$$

Se sugiere que n_j sea tres veces por lo menos la cantidad de parámetros θ_j a estimar y el número puede ser mayor si los grupos no están bien separados.



Cálculo errores de clasificación $f_j, j = 1, 2$

Se definen varios tipos de errores

a) El error actual

$$e_{1,act} = R(\mathcal{P}_1, \hat{G}_0^*) = \int_{\hat{G}_{2,0}} f_1(\mathbf{x}) d\mathbf{x} = 1 - \int_{\hat{G}_{1,0}} f_1(\mathbf{x}) d\mathbf{x}$$

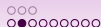
$$e_{2,act} = R(\mathcal{P}_2, \hat{G}_0^*) = \int_{\hat{G}_{1,0}} f_2(\mathbf{x}) d\mathbf{x} = 1 - \int_{\hat{G}_{2,0}} f_2(\mathbf{x}) d\mathbf{x}$$

$$e_{act} = \pi_1 e_{1,act} + \pi_2 e_{2,act}$$

Claramente, $e_{opt} \leq e_{act}$

b) La tasa de error actual esperada

$$\mathbb{E}e_{act} = \pi_1 \mathbb{E}e_{1,act} + \pi_2 \mathbb{E}e_{2,act}$$



Estimación de los errores de clasificación

1) El estimador *plug-in*

$$\hat{e}_{j,act} = \sum_{\ell \neq j} \int_{\hat{G}_{\ell,0}} \hat{f}_j(\mathbf{x}) d\mathbf{x} = 1 - \int_{\hat{G}_{j,0}} \hat{f}_j(\mathbf{x}) d\mathbf{x}$$

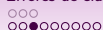
$$\hat{e}_{act} = \sum_{j=1}^k \pi_j \hat{e}_{k,act}$$

Este error se basa en la correcta especificación del modelo pero además en muchos casos, como veremos, subestima el error real e_{opt} .

Si $k = 2$

$$\hat{e}_{1,act} = \int_{\hat{G}_{2,0}} \hat{f}_1(\mathbf{x}) d\mathbf{x} = 1 - \int_{\hat{G}_{1,0}} \hat{f}_1(\mathbf{x}) d\mathbf{x}$$

$$\hat{e}_{2,act} = \int_{\hat{G}_{1,0}} \hat{f}_2(\mathbf{x}) d\mathbf{x} = 1 - \int_{\hat{G}_{2,0}} \hat{f}_2(\mathbf{x}) d\mathbf{x} \quad \hat{e}_{act} = \pi_1 e_{1,act} + \pi_2 e_{2,act}$$



Estimación de los errores de clasificación

2) La tasa de error aparente.

Consideremos la regla basada en las regiones $\hat{G}_{j,0}$, $1 \leq j \leq k$ y sean

$$n_{i,j} = \#\{\mathbf{x}_{i\ell} \text{ clasificadas en la población } \mathcal{P}_j\} = \#\{\mathbf{x}_{i\ell} \in \hat{G}_{j,0}\}$$

$$n_i = \sum_{j=1}^k n_{ij} \text{ el total de observaciones de la población } i\text{-ésima,}$$

$$\hat{\pi}_i = \frac{n_i}{n} \text{ con } n = \sum_{i=1}^k n_i$$

$$m_i = \#\{\mathbf{x}_{i\ell} \text{ mal clasificadas}\} = \sum_{j \neq i} n_{ij}$$



Estimación de los errores de clasificación

2) La tasa de error aparente es

$$e_{i,app} = \frac{m_i}{n_i} \quad e_{app} = \sum_{j=1}^k \pi_j e_{i,app}$$

$$\hat{e}_{app} = \sum_{j=1}^k \hat{\pi}_j e_{i,app} = \frac{\sum_{i=1}^k m_i}{n}$$

El método basado en \hat{e}_{app} se llama también de resustitución. Este estimador del error es muy optimista ya que tiende a subestimar la probabilidad real de error, pues los mismos datos se usan para armar la regla (estimar los parámetros) y para evaluar la regla resultante. Los estimadores de los parámetros obtenidos son los que mejor ajustan a los datos y por ello tiendo a clasificar mejor.



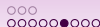
Estimación de los errores de clasificación

- 2) Veamos un ejemplo de como la tasa de error aparente subestima.

Sea $X \in \mathbb{R}$ y supongamos que $n_1 = n_2 = 1$, o sea, tenemos las observaciones x_1 y x_2 de la población 1 y 2, respectivamente. Supongamos $x_1 > x_2$.

Consideremos la regla de clasificación que asigna x a \mathcal{P}_1 si $x \leq (x_1 + x_2)/2$ y al grupo 2 en otro caso.

Entonces, $\hat{e}_{app} = 0$, lo cual es demasiado optimista.



Estimación de los errores de clasificación

3) El estimador de convalidación cruzada.

En este método se sacan las observaciones de a una. Con los $n - 1$ datos restantes se arma la regla y se clasifica la observación extraída. Sea

$$a_i = \#\{\mathbf{x}_{i\ell} \text{ mal clasificadas}, 1 \leq \ell \leq n_i\}$$

$$e_{i,cv} = \frac{a_i}{n_i} \quad e_{cv} = \sum_{i=1}^k \pi_i e_{i,cv}$$

$$\hat{e}_{cv} = \sum_{i=1}^k \hat{\pi}_i e_{i,cv} = \frac{\sum_{i=1}^k a_i}{n}$$



Estimación de los errores de clasificación

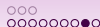
3) El estimador de convalidación cruzada.

Este método da estimadores consistentes del error pero con varianza grande.

Obviamente, es más costoso computacionalmente pero da resultados más honestos y debería ser usado si es posible.

En el caso normal, las fórmulas para los estimadores de los parámetros evitan efectuar el cálculo de la regla en cada paso.

4) Otra opción es usar el método de M -fold que divide la muestra total en M grupos y construye la regla con $M - 1$ grupos mientras clasifica el grupo restante, sucesivamente.



Estimación de los errores de clasificación

5) El estimador bootstrap.

Como $e_{i,app}$ es sesgado, Efron (1979) sugiere estimar su sesgo usando bootstrap.

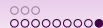
- Para cada $1 \leq i \leq k$, tomamos una muestra $\mathbf{x}_{i\ell}^*$ con reemplazo de la muestra original de la población \mathcal{P}_i , de tamaño n_i .
- Construyamos la regla de clasificación basada en estas muestras que llamaremos $\hat{G}_{i,0}^{*,*}$ con regiones $\hat{G}_{i,0}^*$. Sean

$$m_i^* = \#\{\mathbf{x}_{i\ell}^* \text{ mal clasificadas } 1 \leq \ell \leq n_i\} = \#\{\mathbf{x}_{i\ell}^* \notin \hat{G}_{j,0}^* \mid 1 \leq \ell \leq n_i\}$$

$$m_i^{**} = \#\{\mathbf{x}_{i\ell} \text{ mal clasificadas } 1 \leq \ell \leq n_i\} = \#\{\mathbf{x}_{i\ell} \notin \hat{G}_{j,0}^* \mid 1 \leq \ell \leq n_i\}$$

$$d_i = \frac{m_i^{**} - m_i^*}{n_i}$$

- Repitase a) y b) un número B grande de veces. Sea $d_{i,s}$ el valor de d_i en la replicación s y defina $\bar{d}_i = \sum_{s=1}^B d_{i,s}$



Estimación de los errores de clasificación

5) El estimador bootstrap.

El estimador bootstrap se define como

$$e_{i,boot} = \frac{m_i}{n_i} + \bar{d}_i$$

$$e_{boot} = \sum_{i=1}^k \pi_i e_{i,boot}$$

$$\hat{e}_{boot} = \sum_{i=1}^k \hat{\pi}_i e_{i,boot}$$



$$f_j \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$$

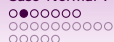
Supongamos que $c_{ij} = 1$ si $i \neq j$ entonces la regla Bayes clasifica $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \mathcal{G}_{i,0}$ donde

$$\begin{aligned} \mathcal{G}_{i,0} &= \{\mathbf{x} \in \mathbb{R}^p : \frac{\pi_\ell}{\pi_i} < \frac{f_i(\mathbf{x})}{f_\ell(\mathbf{x})} \quad \forall \ell \neq i\} \\ &= \{\mathbf{x} \in \mathbb{R}^p : \log\left(\frac{f_i(\mathbf{x})}{f_\ell(\mathbf{x})}\right) > \log\left(\frac{\pi_\ell}{\pi_i}\right) \quad \forall \ell \neq i\} \end{aligned}$$

donde

$$\begin{aligned} \log\left(\frac{f_i(\mathbf{x})}{f_\ell(\mathbf{x})}\right) &= \frac{1}{2} \log \frac{\det(\boldsymbol{\Sigma}_\ell)}{\det(\boldsymbol{\Sigma}_j)} - \frac{1}{2} (\boldsymbol{\mu}_i^T \boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \boldsymbol{\mu}_\ell^T \boldsymbol{\Sigma}_\ell^{-1} \boldsymbol{\mu}_\ell) \\ &\quad - \frac{1}{2} \{\mathbf{x}^T (\boldsymbol{\Sigma}_i^{-1} - \boldsymbol{\Sigma}_\ell^{-1}) \mathbf{x} - 2\mathbf{x}^T (\boldsymbol{\Sigma}_i^{-1} \boldsymbol{\mu}_i - \boldsymbol{\Sigma}_\ell^{-1} \boldsymbol{\mu}_\ell)\} \end{aligned}$$

o sea, obtenemos una forma cuadrática en \mathbf{x} .



Caso $\Sigma_\ell = \Sigma$

En este caso,

$$\log \left(\frac{f_i(\mathbf{x})}{f_\ell(\mathbf{x})} \right) = \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i) - \boldsymbol{\mu}_\ell^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_\ell)$$

Por lo tanto, si llamamos

$$L_i(\mathbf{x}) = \log \pi_i + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i)$$

asigno \mathbf{x} al grupo con mayor $L_i(\mathbf{x})$, o sea, clasifico $\mathbf{x} \in \mathcal{P}_i$ si

$$L_i(\mathbf{x}) = \max_{1 \leq \ell \leq k} L_\ell(\mathbf{x})$$

Si $\pi_j = \frac{1}{k}$ para todo j , esta regla de clasificación es la obtenida antes con las coordenadas discriminantes.



Caso $\Sigma_\ell = \Sigma$

En este caso, asigno \mathbf{x} al grupo con mayor $L_i(\mathbf{x})$ donde llamamos

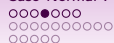
$$L_i(\mathbf{x}) = \log \pi_i + \boldsymbol{\mu}_i^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_i \right)$$

Las funciones

$$d_{i\ell}(\mathbf{x}) = L_i(\mathbf{x}) - L_\ell(\mathbf{x}) = (\boldsymbol{\mu}_i - \boldsymbol{\mu}_\ell)^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \frac{(\boldsymbol{\mu}_i + \boldsymbol{\mu}_\ell)}{2} \right) + \log \pi_i - \log \pi_\ell$$

se llaman funciones discriminantes y $d_{i\ell}(\mathbf{x}) = -d_{\ell i}(\mathbf{x})$. Sea $\boldsymbol{\alpha}_{i,\ell} = \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_i - \boldsymbol{\mu}_\ell)$, Luego

$$d_{i\ell}(\mathbf{x}) = \boldsymbol{\alpha}_{i,\ell}^T \left(\mathbf{x} - \frac{(\boldsymbol{\mu}_i + \boldsymbol{\mu}_\ell)}{2} \right) + \log \pi_i - \log \pi_\ell$$



Caso $\Sigma_\ell = \Sigma$

Si tenemos k poblaciones, sólo necesitamos encontrar

$r = \min(p, k - 1)$ direcciones de proyección en lugar de $\begin{pmatrix} k \\ 2 \end{pmatrix}$.

Efectivamente, basta conocer $\alpha_{i,i+1}$, $1 \leq i \leq k - 1$ ya que

$$\alpha_{i,i+2} = \alpha_{i,i+1} - \alpha_{i+1,i+2}$$

Ejemplo: Si $k = 3$ y obtenemos que $L_1(\mathbf{x}) > L_2(\mathbf{x})$ y $L_2(\mathbf{x}) > L_3(\mathbf{x})$ entonces $L_1(\mathbf{x}) > L_3(\mathbf{x})$.

Si además $p = 2$ cada ecuación $d_{i\ell}(\mathbf{x}) = 0$ es una recta y las tres rectas se cortan en el mismo punto ya que

$$d_{13}(\mathbf{x}) = L_1(\mathbf{x}) - L_3(\mathbf{x}) = L_1(\mathbf{x}) - L_2(\mathbf{x}) + (L_2(\mathbf{x}) - L_3(\mathbf{x})) = d_{12}(\mathbf{x}) + d_{23}(\mathbf{x})$$

○○○○○
 ○○○○○○○○○○○○
 ○○○○○○○○○○

○○○○○
 ○○○○

○○○
 ○○○○○○○○

○○○○●○○
 ○○○○○○○○
 ○○○○

○○
 ○○○○○○
 ○○

Caso $\Sigma_\ell = \Sigma$

Si $k = 2$, como $\alpha = \Sigma^{-1}(\mu_1 - \mu_2)$

$$d_{12}(\mathbf{x}) = \alpha^T \left(\mathbf{x} - \frac{(\mu_1 + \mu_2)}{2} \right) + \log \pi_1 - \log \pi_2$$

es la regla discriminante lineal de Fisher que clasifica en el grupo 1 si $d_{12}(\mathbf{x}) > 0$.

El hiperplano $d_{12}(\mathbf{x}) = 0$ determina un hiperplano que separa los dos grupos.



Caso $\Sigma_\ell = \Sigma$

Veamos que si $\pi_j = \frac{1}{k}$ para todo j , la regla de clasificación Bayes es la obtenida antes con las coordenadas discriminantes. Sea

- $\Sigma_B = \sum_{i=1}^k \pi_i (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T$, $s = \text{rango}(\Sigma_B)$
- $\mathbf{z} = \mathbf{z} = \mathbf{A}^T \mathbf{x} = (\mathbf{z}_1^T, \mathbf{z}_2^T)^T$ el vector de variables discriminantes con

$$\mathbf{A} = (\alpha_1, \dots, \alpha_p) = (\mathbf{A}_1, \mathbf{A}_2) \quad \text{donde} \quad \mathbf{A}_1 = (\alpha_1, \dots, \alpha_s)$$

- $\mathbf{z}^{(1)} = \mathbf{A}_1^T \mathbf{x}$, $\mathbf{z}^{(2)} = \mathbf{A}_2^T \mathbf{x}$
- $\nu_i = \mathbf{A}^T \mu_i$, entonces $\nu_i^{(2)} = \nu_i^{(2)}$

Si $\mathbf{x}|G = i \sim N(\mu_i, \Sigma)$ vemos que

$$\mathbf{z}|G = i \sim N(\nu_i, \mathbf{I}_p)$$

o sea,

$$\mathbf{z}^{(1)}|G = i \sim N(\nu_i^{(1)}, \mathbf{I}_s) \quad \mathbf{z}^{(2)} \sim N(\nu_i^{(2)}, \mathbf{I}_{p-s})$$



Caso $\Sigma_\ell = \Sigma$

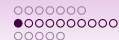
La regla que vimos asignaba \mathbf{x}_0 al grupo i si $\mathbf{v}_0 = \mathbf{A}_1^T \mathbf{x}_0 \in \mathcal{G}_i$ donde

$$\begin{aligned} \mathcal{G}_i &= \{ \mathbf{v} \in \mathbb{R}^s : \|\mathbf{v} - \boldsymbol{\nu}_i^{(1)}\| < \|\mathbf{v} - \boldsymbol{\nu}_\ell^{(1)}\| \quad \forall \ell \neq i \} \\ &= \{ \mathbf{v} \in \mathbb{R}^s : (\boldsymbol{\nu}_i^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_i^{(1)}) > (\boldsymbol{\nu}_\ell^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_\ell^{(1)}) \quad \forall \ell \neq i \} \end{aligned}$$

Esta regla es Bayes cuando $\pi_j = \frac{1}{k}$ para todo j .

Para una probabilidad a priori τ general tenemos que modificar \mathcal{G}_i por $\mathcal{G}_{i,\tau}$, o sea, asigno \mathbf{x}_0 al grupo i si $\mathbf{v}_0 = \mathbf{A}_1^T \mathbf{x}_0 \in \mathcal{G}_{i,\tau}$ donde

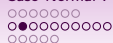
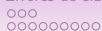
$$\begin{aligned} \mathcal{G}_{i,\tau} &= \{ \mathbf{v} \in \mathbb{R}^s : \|\mathbf{v} - \boldsymbol{\nu}_i^{(1)}\|^2 - 2 \log(\pi_i) < \|\mathbf{v} - \boldsymbol{\nu}_\ell^{(1)}\|^2 - 2 \log(\pi_\ell) \quad \forall \ell \neq i \} \\ &= \{ \mathbf{v} \in \mathbb{R}^s : (\boldsymbol{\nu}_i^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_i^{(1)}) > (\boldsymbol{\nu}_\ell^{(1)})^T (\mathbf{v} - \frac{1}{2} \boldsymbol{\nu}_\ell^{(1)}) \\ &\quad + \log(\pi_\ell) - \log(\pi_i) \quad \forall \ell \neq i \} \end{aligned}$$



$$f_j \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j), j = 1, 2$$

Definición. Supongamos que $\mathbf{x} \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, $j = 1, 2$. Se dice que \mathbf{x} está dada en forma canónica si $\boldsymbol{\mu}_1 = \mathbf{0}$, $\boldsymbol{\Sigma}_1 = \mathbf{I}_p$ y $\boldsymbol{\Sigma}_2$ es diagonal. Llamaremos $\boldsymbol{\Lambda} = \boldsymbol{\Sigma}_2 = \text{diag}(\lambda_1, \dots, \lambda_p)$ y $\boldsymbol{\mu}_2 = \boldsymbol{\delta} = (\delta_1, \dots, \delta_p)^T$.

Propiedad. Supongamos que $\mathbf{x} \sim N(\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)$, $j = 1, 2$ donde $\boldsymbol{\Sigma}_1 > 0$, $\boldsymbol{\Sigma}_2 > 0$. Entonces existe $\boldsymbol{\Gamma} \in \mathbb{R}^{p \times p}$ no singular tal que $\mathbf{z} = \boldsymbol{\Gamma}^T(\mathbf{x} - \boldsymbol{\mu}_1)$ está en forma canónica.



Por lo tanto, si $c_{ij} = 1$ si $i \neq j$ entonces la regla Bayes clasifica $\mathbf{x} \in \mathcal{P}_1$ si $\mathbf{z} = \mathbf{\Gamma}^T(\mathbf{x} - \boldsymbol{\mu}_1) \in \mathcal{G}_{1,0}$ donde

$$\mathcal{G}_{1,0} = \{\mathbf{z} \in \mathbb{R}^p : \log \left(\frac{f_1(\mathbf{x})}{f_2(\mathbf{x})} \right) > \log \left(\frac{\pi_2}{\pi_1} \right)\} = \{\mathbf{z} \in \mathbb{R}^p : Q(\mathbf{z}) > 0\}$$

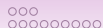
con

$$Q(\mathbf{z}) = \sum_{i=1}^p a_{ii} z_i^2 + \sum_{i=1}^p b_i z_i + c$$

$$a_{ii} = \frac{1}{2} \left(\frac{1}{\lambda_i} - 1 \right) \quad b_i = - \frac{\delta_i}{\lambda_i}$$

$$c = \log \left(\frac{\pi_1}{\pi_2} \right) + \frac{1}{2} \sum_{s=1}^p \log \lambda_s + \frac{1}{2} \sum_{s=1}^p \frac{\delta_s^2}{\lambda_s}$$

La ventaja de la forma canónica es que los términos de la forma $a_{ij} x_i x_j$ desaparecen lo que hace más fácil de entender la regla de clasificación.



Ejemplo $p = 2$

Tomemos $\pi_1 = \pi_2 = 1/2$ y los siguientes valores para δ y Λ

Caso	δ_1	δ_2	λ_1	λ_2	a_{11}	a_{22}	b_1	b_2	c
A	2	1	0.4	0.1	0.75	4.5	-5	-10	8.391
B	2	-1	4	0.25	-0.375	1.5	-0.5	4	2.5
C	3	1	4	1	-0.375	0	-0.75	-1	2.318
D	2	0	10	1	-0.45	0	-0.2	0	1.351

```

○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○
○○○○○

```

```

○○○
○○○○○○○○

```

```

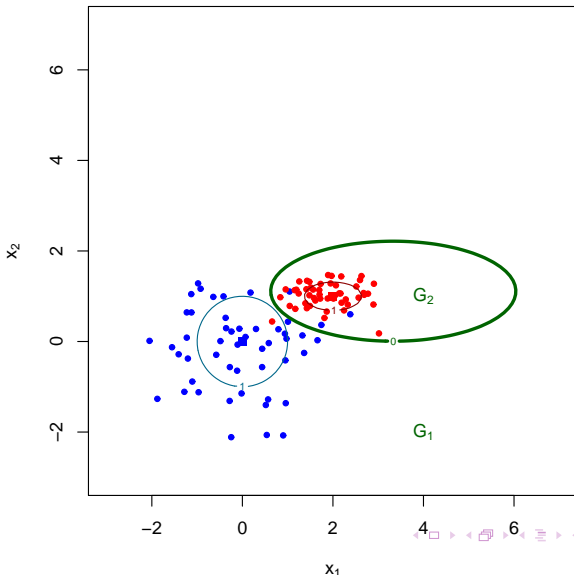
○○○○○○
○○●○○○○○
○○○○○

```

```

○○
○○○○○○○
○○○

```

Ejemplo $p = 2$, Caso A

```

○○○○○○
○○○○○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○○
○○○○○○

```

```

○○○
○○○○○○○○○○

```

```

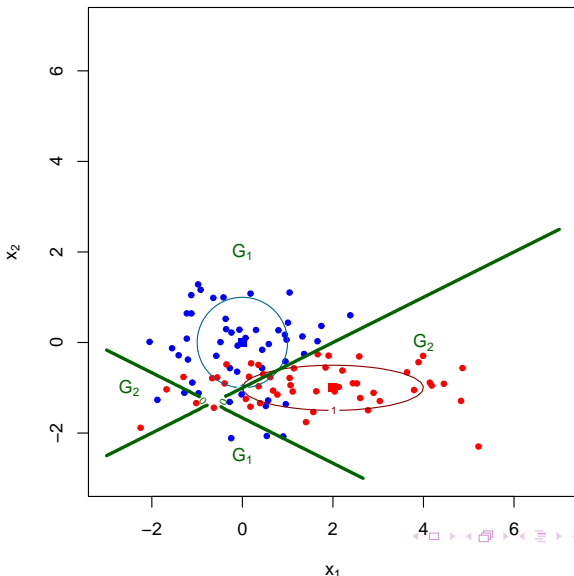
○○○○○○○
○○○○●○○○○
○○○○○

```

```

○○
○○○○○○○○
○○○

```

Ejemplo $p = 2$, Caso B

```

○○○○○○
○○○○○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○○
○○○○○○

```

```

○○○
○○○○○○○○○○

```

```

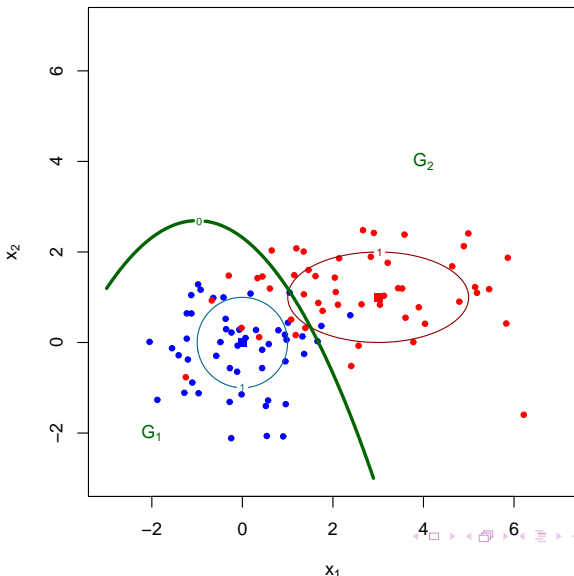
○○○○○○○
○○○○○○●○○○○
○○○○○

```

```

○○
○○○○○○○○
○○○

```

Ejemplo $p = 2$, Caso C, $a_{22} = 0$ 

```

○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○
○○○○○

```

```

○○○
○○○○○○○○○

```

```

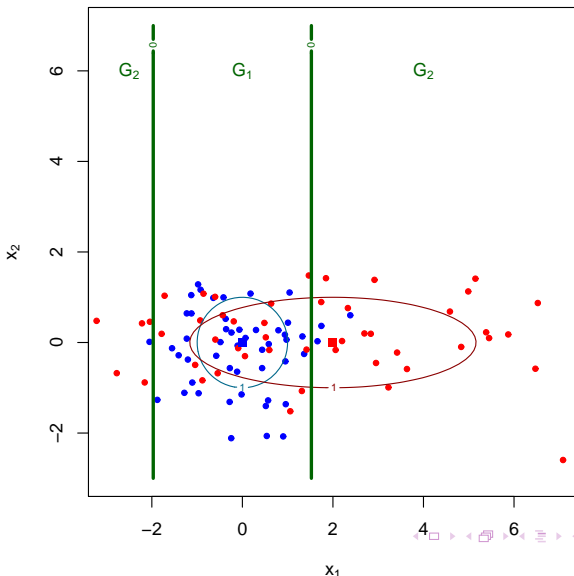
○○○○○○○
○○○○○○●○○○
○○○○○

```

```

○○
○○○○○○○
○○○

```

Ejemplo $p = 2$, Caso D, $a_{22} = 0$, $b_2 = 0$ 

```

○○○○○○
○○○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○
○○○○○

```

```

○○○
○○○○○○○○○

```

```

○○○○○○○
○○○○○○●○○
○○○○○

```

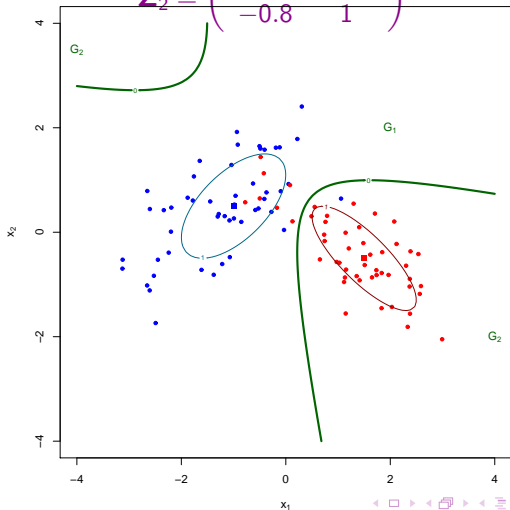
```

○○
○○○○○○○
○○○

```

Ejemplo $\mu_1 = (-1, 0.5)$, $\mu_2 = (1.5, -0.5)$, $\Sigma_1 = \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix}$,

$$\Sigma_2 = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}$$



```

○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○
○○○○○

```

```

○○○
○○○○○○○○○

```

```

○○○○○○○
○○○○○○○○●○
○○○○○

```

```

○○
○○○○○○○
○○○

```

Ejemplo $\mu_1 = (-1, 0.5)$, $\mu_2 = (1.5, -0.5)$, $\Sigma_1 = \begin{pmatrix} 1 & 0.6 \\ 0.6 & 1 \end{pmatrix}$,

$$\Sigma_2 = \begin{pmatrix} 1 & -0.8 \\ -0.8 & 1 \end{pmatrix}$$

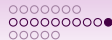
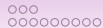
Datos transformados $\mathbf{z} = \Gamma^T(\mathbf{x} - \mu_1)$.

$$\Gamma = \begin{pmatrix} -0.5590 & -1.1180 \\ -0.5590 & 1.1180 \end{pmatrix}$$

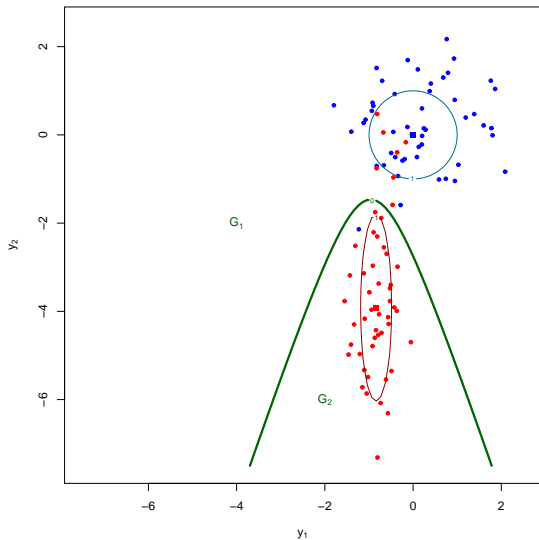
$$\lambda_1 = 0.125 \quad \lambda_2 = 4.5$$

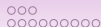
$$a_{11} = 3.500 \quad a_{22} = -0.3889$$

$$b_1 = 6.708 \quad b_2 = 0.8696 \quad c = 5.325$$



Datos transformados $\mathbf{z} = \mathbf{\Gamma}^T(\mathbf{x} - \boldsymbol{\mu}_1)$





Cálculo errores de clasificación $f_j \sim N_p(\mu_j, \Sigma)$, $j = 1, 2$

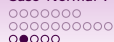
- $\alpha = \Sigma^{-1}(\mu_1 - \mu_2)$,
- la regla G_0^* clasifica en el grupo 1 si $\mathbf{x} \in \mathcal{G}_{1,0}$, con $\mathcal{G}_{1,0} = \{\mathbf{x} : d_{12}(\mathbf{x}) > 0\}$

$$d_{12}(\mathbf{x}) = \alpha^T \left(\mathbf{x} - \frac{(\mu_1 + \mu_2)}{2} \right) + \log \pi_1 - \log \pi_2$$

$$\Delta_p^2 = \alpha^T (\mu_1 - \mu_2) = (\mu_1 - \mu_2)^T \Sigma^{-1} (\mu_1 - \mu_2) = \alpha^T \Sigma \alpha$$

Luego

$$R(\mathcal{P}_1, G_0^*) = \int_{\mathcal{G}_{2,0}} f_1(\mathbf{x}) d\mathbf{x} = 1 - \int_{\mathcal{G}_{1,0}} f_1(\mathbf{x}) d\mathbf{x}$$



Cálculo errores de clasificación $f_j \sim N_p(\mu_j, \Sigma)$, $j = 1, 2$

En \mathcal{P}_1 , $\mathbf{x} \sim N_p(\mu_1, \Sigma)$ luego

$$\alpha^T \mathbf{x} \sim N(\alpha^T \mu_1, \alpha^T \Sigma \alpha) = N(\alpha^T \mu_1, \Delta_p^2)$$

de donde

$$R(\mathcal{P}_1, G_0^*) = \Phi \left(\frac{\log \left(\frac{\pi_2}{\pi_1} \right) - \frac{1}{2} \Delta_p^2}{\Delta_p} \right)$$

$$R(\mathcal{P}_2, G_0^*) = \Phi \left(- \frac{\log \left(\frac{\pi_2}{\pi_1} \right) + \frac{1}{2} \Delta_p^2}{\Delta_p} \right)$$

$$e_{opt} = r(\tau, G_0^*) = \pi_1 \Phi \left(\frac{\log \left(\frac{\pi_2}{\pi_1} \right) - \frac{1}{2} \Delta_p^2}{\Delta_p} \right) + \pi_2 \Phi \left(- \frac{\log \left(\frac{\pi_2}{\pi_1} \right) + \frac{1}{2} \Delta_p^2}{\Delta_p} \right)$$



Cálculo errores de clasificación $f_j \sim N_p(\mu_j, \Sigma)$, $j = 1, 2$

Si $\pi_1 = \pi_2 = 1/2$ entonces

$$R(\mathcal{P}_1, G_0^*) = R(\mathcal{P}_2, G_0^*) = \Phi\left(-\frac{1}{2}\Delta_p\right)$$

y G_0^* es minimax.

En general, la regla minimax asigna al grupo 1 si

$D(\mathbf{x}) = \alpha^T \left(\mathbf{x} - \frac{(\mu_1 + \mu_2)}{2}\right) > \log(c)$ donde c se elige de modo que

$$\Phi\left(\frac{\log(c) - \frac{1}{2}\Delta_p^2}{\Delta_p}\right) = \Phi\left(\frac{-\log(c) - \frac{1}{2}\Delta_p^2}{\Delta_p}\right)$$

que tiene como solution $c = 1$ coincidiendo con el método de cociente de verosimilitud.



Cálculo errores de clasificación $f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$, $j = 1, \dots, k$

- la regla G_0^* clasifica en el grupo i si $\mathbf{x} \in \mathcal{G}_{i,0}$, con

$$\mathcal{G}_{i,0} = \{\mathbf{x} : d_{i\ell}(\mathbf{x}) > 0 \forall \ell \neq i\} = \{\mathbf{x} : L_i(\mathbf{x}) = \max_{\ell} L_{\ell}(\mathbf{x})\}$$

$$d_{i\ell}(\mathbf{x}) = L_i(\mathbf{x}) - L_{\ell}(\mathbf{x}) = (\boldsymbol{\mu}_i - \boldsymbol{\mu}_{\ell})^T \boldsymbol{\Sigma}^{-1} \left(\mathbf{x} - \frac{(\boldsymbol{\mu}_i + \boldsymbol{\mu}_{\ell})}{2} \right) + \log \pi_i - \log \pi_{\ell}$$

- $\Delta_{i\ell}^2 = (\boldsymbol{\mu}_i - \boldsymbol{\mu}_{\ell})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_{\ell})$

Luego

$$R(\mathcal{P}_i, G_0^*) = \sum_{\ell \neq i} \int_{\mathcal{G}_{\ell,0}} f_i(\mathbf{x}) d\mathbf{x} = 1 - \int_{\mathcal{G}_{i,0}} f_i(\mathbf{x}) d\mathbf{x}$$



Cálculo errores de clasificación $f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$, $j = 1, \dots, k$

En \mathcal{P}_i , $\mathbf{x} \sim N_p(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ luego

$$d_{i\ell}(\mathbf{x}) \sim N\left(\frac{1}{2}\Delta_{i\ell}^2 + \log \pi_i - \log \pi_\ell, \Delta_{i\ell}^2\right)$$

Más aún, el vector $\mathbf{d}_i(\mathbf{x}) = (d_{i\ell}(\mathbf{x}))_{\ell \neq i}$ tiene distribución normal $(k-1)$ -variada y

$$\text{COV}(d_{i\ell}(\mathbf{x}), d_{ij}(\mathbf{x})) = (\boldsymbol{\mu}_i - \boldsymbol{\mu}_\ell)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_i - \boldsymbol{\mu}_j)$$

Si $k = 3$ se pueden calcular, fácilmente.



$f_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ con $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i$ desconocidos

Supongamos tener $\mathbf{x}_{ij} \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$, $1 \leq j \leq n_i$, $1 \leq i \leq k$ entonces estimamos $\boldsymbol{\mu}_i$ y $\boldsymbol{\Sigma}_i$ por

$$\hat{\boldsymbol{\mu}}_i = \bar{\mathbf{x}}_i \quad \mathbf{S}_i = \frac{\mathbf{Q}_i}{n_i - 1} = \frac{1}{n_i} \sum_{j=1}^{n_i} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)(\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)^T$$

Clasificamos $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \hat{\mathcal{G}}_{i,0}$ con

$$\hat{\mathcal{G}}_{i,0} = \left\{ \mathbf{x} \in \mathbb{R}^p : \log \left(\frac{\hat{f}_i(\mathbf{x})}{\hat{f}_\ell(\mathbf{x})} \right) > \log \left(\frac{\pi_\ell}{\pi_i} \right) \quad \forall \ell \neq i \right\}$$

donde $\hat{f}_i(\mathbf{x}) = f(\mathbf{x}, \hat{\boldsymbol{\mu}}_i, \mathbf{S}_i)$ con $f(\mathbf{x}, \boldsymbol{\mu}, \boldsymbol{\Sigma}) \sim N(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

$$\begin{aligned} \log \left(\frac{\hat{f}_i(\mathbf{x})}{\hat{f}_\ell(\mathbf{x})} \right) &= \frac{1}{2} \log \frac{\det(\mathbf{S}_\ell)}{\det(\mathbf{S}_i)} - \frac{1}{2} (\hat{\boldsymbol{\mu}}_i^T \mathbf{S}_i^{-1} \hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\mu}}_\ell^T \mathbf{S}_\ell^{-1} \hat{\boldsymbol{\mu}}_\ell) \\ &\quad - \frac{1}{2} \left\{ \mathbf{x}^T (\mathbf{S}_i^{-1} - \mathbf{S}_\ell^{-1}) \mathbf{x} - 2\mathbf{x}^T (\mathbf{S}_i^{-1} \hat{\boldsymbol{\mu}}_i - \mathbf{S}_\ell^{-1} \hat{\boldsymbol{\mu}}_\ell) \right\} \end{aligned}$$



$f_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$ con $\boldsymbol{\mu}_i, \boldsymbol{\Sigma}$ desconocidos

Supongamos tener $\mathbf{x}_{ij} \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$, $1 \leq j \leq n_i$, $1 \leq i \leq k$ entonces estimamos $\boldsymbol{\mu}_i$ y $\boldsymbol{\Sigma}$ por

$$\hat{\boldsymbol{\mu}}_i = \bar{\mathbf{x}}_i \quad \mathbf{S} = \frac{1}{n-k} \sum_{i=1}^k \mathbf{Q}_i = \frac{1}{n-k} \sum_{i=1}^k \sum_{j=1}^{n_i} (\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)(\mathbf{x}_{i,j} - \bar{\mathbf{x}}_i)^T$$

Clasificamos $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \hat{\mathcal{G}}_{i,0}$ con

$$\hat{\mathcal{G}}_{i,0} = \{\mathbf{x} \in \mathbb{R}^p : \hat{L}_i(\mathbf{x}) > \hat{L}_\ell(\mathbf{x}) \quad \forall \ell \neq i\} = \{\mathbf{x} \in \mathbb{R}^p : \hat{d}_{i\ell}(\mathbf{x}) > 0\}$$

donde

$$\hat{L}_i(\mathbf{x}) = \log \pi_i + \hat{\boldsymbol{\mu}}_i^T \mathbf{S}^{-1} \left(\mathbf{x} - \frac{1}{2} \hat{\boldsymbol{\mu}}_i \right)$$

$$\begin{aligned} \hat{d}_{i\ell}(\mathbf{x}) &= \hat{L}_i(\mathbf{x}) - \hat{L}_\ell(\mathbf{x}) \\ &= \log \left(\frac{\pi_i}{\pi_\ell} \right) + (\hat{\boldsymbol{\mu}}_i - \hat{\boldsymbol{\mu}}_\ell)^T \mathbf{S}^{-1} \left(\mathbf{x} - \frac{\hat{\boldsymbol{\mu}}_i + \hat{\boldsymbol{\mu}}_\ell}{2} \right) \end{aligned}$$



Estimación de los errores de clasificación $k = 2$, $f_j \sim N_p(\mu_j, \Sigma)$

En el caso normal clasificamos $\mathbf{x} \in \mathcal{P}_1$ si $\hat{d}_{12}(\mathbf{x}) > 0$,

$$\begin{aligned}\hat{d}_{12}(\mathbf{x}) &= \log\left(\frac{\pi_1}{\pi_2}\right) + (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}^{-1} \left(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}\right) \\ &= \log\left(\frac{\pi_1}{\pi_2}\right) + \hat{\boldsymbol{\alpha}}^T \left(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}\right)\end{aligned}$$

Por lo tanto, como si $\mathbf{x} \sim N_p(\mu_1, \Sigma)$ y si llamamos $\sigma^2 = \hat{\boldsymbol{\alpha}}^T \Sigma \hat{\boldsymbol{\alpha}}$, $\mathbf{X}_1 = (\mathbf{x}_{1,1}, \dots, \mathbf{x}_{1,n_1})$ y $\mathbf{X}_2 = (\mathbf{x}_{2,1}, \dots, \mathbf{x}_{2,n_2})$, entonces

$$\begin{aligned}\hat{\boldsymbol{\alpha}}^T \left(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}\right) | (\mathbf{X}_1, \mathbf{X}_2) &\sim N\left(\hat{\boldsymbol{\alpha}}^T \left(\mu_1 - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2}\right), \sigma^2\right) \\ \hat{d}_{12}(\mathbf{x}) &\sim N\left(\hat{d}_{12}(\mu_1), \sigma^2\right)\end{aligned}$$



Estimación de los errores de clasificación $k = 2$,

$$f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$$

$$\begin{aligned}
 e_{1,act} &= \int_{\hat{\mathcal{G}}_{2,0}} f_1(\mathbf{x}) d\mathbf{x} = \int_{\hat{d}_{12}(\mathbf{x}) < 0} f_1(\mathbf{x}) d\mathbf{x} \\
 &= \mathbb{P} \left(\hat{\boldsymbol{\alpha}}^T \left(\mathbf{x} - \frac{\bar{\mathbf{x}}_1 + \bar{\mathbf{x}}_2}{2} \right) + \log \left(\frac{\pi_1}{\pi_2} \right) < 0 \mid \mathbf{x} \sim N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}) \right) \\
 &= \Phi \left(-\frac{\hat{d}_{12}(\boldsymbol{\mu}_1)}{\sigma} \right) = \Phi \left(\frac{\hat{d}_{21}(\boldsymbol{\mu}_1)}{\sigma} \right)
 \end{aligned}$$

Por lo tanto,

$$e_{act} = \pi_1 \Phi \left(\frac{\hat{d}_{21}(\boldsymbol{\mu}_1)}{\sigma} \right) + \pi_2 \Phi \left(\frac{\hat{d}_{12}(\boldsymbol{\mu}_2)}{\sigma} \right)$$



Estimación de los errores de clasificación $k = 2$,
 $f_j \sim N_p(\mu_j, \Sigma)$

$$\hat{e}_{act} = \pi_1 \Phi \left(\frac{\log \left(\frac{\pi_2}{\pi_1} \right) - \frac{1}{2} D_p^2}{D_p} \right) + \pi_2 \Phi \left(- \frac{\log \left(\frac{\pi_2}{\pi_1} \right) + \frac{1}{2} D_p^2}{D_p} \right)$$

donde

$$D_p^2 = \hat{\alpha}^T (\hat{\mu}_1 - \hat{\mu}_2) = (\hat{\mu}_1 - \hat{\mu}_2)^T \mathbf{S}^{-1} (\hat{\mu}_1 - \hat{\mu}_2)$$

Si π_1 y π_2 son desconocidos se estiman por n_1/n y n_2/n , con $n = n_1 + n_2$. Si $\pi_1 = \pi_2 = 0.5$, entonces

$$\hat{e}_{act} = \Phi \left(- \frac{1}{2} D_p \right)$$



Estimación de los errores de clasificación $k = 2$,

$$f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$$

Como $\mathbb{E}\mathcal{F}_{\nu_1, \nu_2}(\lambda) = \nu_2(\nu_1 + \lambda) / [\nu_1(\nu_2 - 2)]$ si $\nu_2 > 2$ y

$$\frac{n_1 + n_2 - p - 1}{(n_1 + n_2 - 2)p} T_0^2 \sim \mathcal{F}_{p, n_1 + n_2 - p - 1}(\lambda^2)$$

con $\lambda^2 = \frac{n_1 n_2}{n_1 + n_2} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$ y

$$T_0^2 = \frac{n_1 n_2}{n_1 + n_2} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)^T \mathbf{S}^{-1} (\bar{\mathbf{x}}_1 - \bar{\mathbf{x}}_2)$$

tenemos que

$$\mathbb{E}D_p^2 = \frac{n - 2}{n - p - 3} \left(\Delta_p^2 + \frac{pn}{n_1 n_2} \right)$$

D_p^2 sobreestima a Δ_p^2 y por lo tanto, \hat{e}_{act} subestima el error real

e_{act} .



Estimación de los errores de clasificación $k = 2$, $\mathbf{x} | \mathbf{x} \in \mathcal{P}_j \sim f_j = f_j(\mathbf{x}, \boldsymbol{\theta}_j)$

Supongamos que $c_{ij} = 1$ si $i \neq j$ entonces la regla Bayes clasificaba $\mathbf{x} \in \mathcal{P}_i$ si $\mathbf{x} \in \mathcal{G}_{i,0}$ donde

$$\begin{aligned} \mathcal{G}_{i,0} &= \{\mathbf{x} \in \mathbb{R}^P : \pi_\ell f_\ell(\mathbf{x}) < \pi_i f_i(\mathbf{x}) \quad \forall \ell \neq i\} \\ &= \{\mathbf{x} \in \mathbb{R}^P : r_\ell(\mathbf{x}) < r_i(\mathbf{x}) \quad \forall \ell \neq i\} \end{aligned}$$

es decir, clasifico $\mathbf{x} \in \mathcal{P}_i$ si $r_i(\mathbf{x}) = \max_{1 \leq \ell \leq k} r_\ell(\mathbf{x})$, donde $r_i(\mathbf{x}) = f_i(\mathbf{x}, \boldsymbol{\theta}_i)\pi_i$.

Lema. Sea $\hat{r}_i(\mathbf{x}) = f_i(\mathbf{x}, \hat{\boldsymbol{\theta}}_i)\pi_i$ un estimador insesgado de $r_i(\mathbf{x})$, o sea, $\mathbb{E}_{\hat{\boldsymbol{\theta}}} \hat{r}_i(\mathbf{x}) | \mathbf{x} = r_i(\mathbf{x})$ para casi todo \mathbf{x} . Luego

$$\mathbb{E} \hat{e}_{act} \leq e_{opt} < \mathbb{E} e_{act}$$

```

○○○○○○
○○○○○○○○○○○○○○
○○○○○○○○○○

```

```

○○○○○
○○○○○

```

```

○○○
○○○○○○○○

```

```

○○○○○○
○○○○○○○○
○○○○○○○○

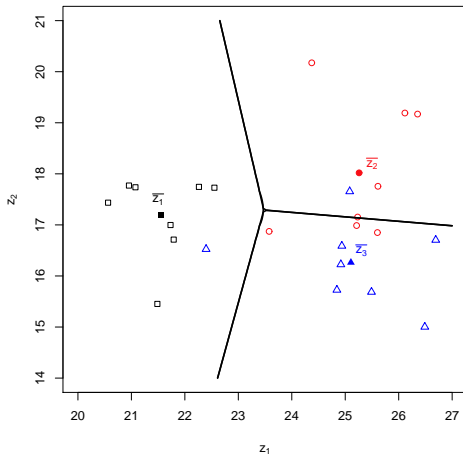
```

```

○○
○○○○○●○○
○○○

```

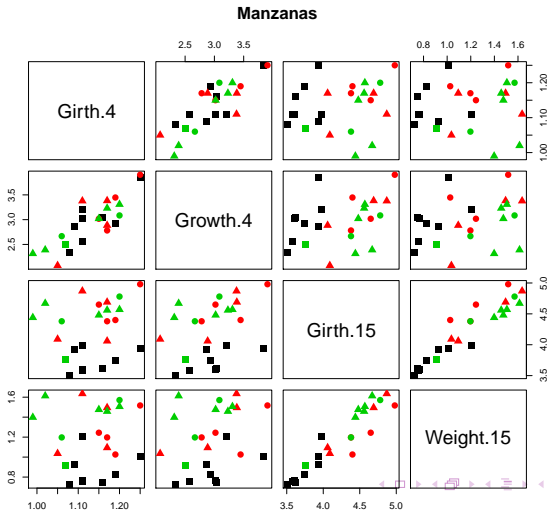
Ejemplo Manzanos





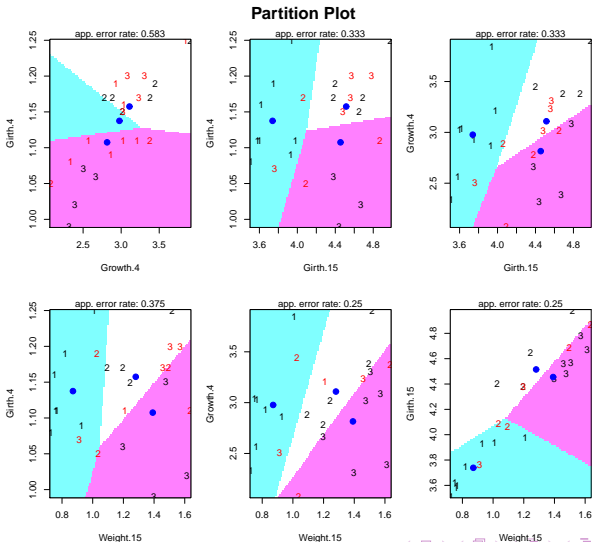
Ejemplo Manzanas

■: Asignado al grupo 1, ● asignado al grupo 2, ▲ asignado al grupo 3





Ejemplo Manzanos: En negro bien clasificados, en rojo mal clasificados



○○○○○
 ○○○○○○○○○○○○
 ○○○○○○○○○○

○○○○○
 ○○○○

○○○
 ○○○○○○○○

○○○○○○○
 ○○○○○○○○
 ○○○○

○○
 ○○○○○○
 ●○○

Comparación entre LDF y QDF, $k = 2$

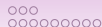
- En general la decisión de elegir entre la regla lineal (LDF) y cuadrática (QDF) se hace en base al resultado del test para $H_0 : \Sigma_1 = \Sigma_2$. Si el test rechaza se usa QDF.
- A pesar de que esta decisión es razonable ya que LDF es óptima si H_0 es cierta, hay un número importante de trabajos que muestran que aunque no lo sea, LDF es tan buena como QDF.
- Uno podría basar su decisión en elegir el método que da menor error aparente, \hat{e}_{app} , lo cual es peligroso ya que este estimador del error subestima el error real, e_{act} . El cálculo del error actual esperado sólo puede hacerse por simulación.
- Una opción es utilizar el e_{cv} para elegir entre ambas reglas.



Comparación entre LDF y QDF, $k = 2$

En general, LDF es buena para pequeños alejamientos de $H_0 : \Sigma_1 = \Sigma_2$. El mejor comportamiento de QDF depende del tamaño de las muestras y de la dimensión.

- Para n_1 y n_2 pequeñas y $p \leq 6$ hay poca pérdida al elegir LDF.
- Para $n_1, n_2 \leq 25$ y p grande y/o diferencias entre Σ_1 y Σ_2 , LDF es preferible.
- Sin embargo, cuando p grande y Σ_1 y Σ_2 son muy distintas, las probabilidades de mala clasificación $e_{1,act}$ y $e_{2,act}$ pueden ser muy grandes para un uso práctico



Comparación entre LDF y QDF, $k = 2$

- Si $\Sigma_1 \neq \Sigma_2$ y la diferencia es grande y $p > 6$, QDF es mucho mejor que LDF si el tamaño de muestra es grande.
 - Se recomienda para $p = 4$, $n_1 = n_2 = 25$
 - 25 observaciones adicionales cada dos dimensiones, o sea, para $p = 6, 8, 10$ se necesitan $n_1 = n_2 = 50, 75, 100$
- Para $n_i \geq 100$ y p moderado los resultados asintóticos que favorecen QDF se alcanzan bastante rápido.
- QDF se deteriora rápidamente si p crece porque S_i no provee una estimación confiable de Σ_i si p es una fracción moderada de n_i .