



# Discriminación

**Graciela Boente**



## 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

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

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

## Definiciones

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

Definimos

a) La función de pérdida como

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

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

donde  $c_{jj} = 0$ .

b) El riesgo de  $G^*$  es

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

## Definiciones

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

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

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

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

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  $\tau$  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  $\tau$  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  $\sum_{j=1}^k \pi_j c_{i|j} h_j(\mathbf{x}_1) = \min_{\ell} \sum_{j=1}^k \pi_j c_{\ell|j} h_j(\mathbf{x}_1)$ .

## 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  $j$

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

Si  $k = 2$ , llamaremos

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

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

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

## 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

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

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

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

## $f_j = f_j(\cdot, \theta_j)$ , $\theta_j$ 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  $\theta_j$  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  $c_{i|j} = 1$  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  $\theta_j$  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  $k = 2$

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

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

## Estimación de los errores de clasificación

Sea

$$\hat{r}_i(\mathbf{x}) = \pi_i f(\mathbf{x}, \hat{\theta}_i)$$

un estimador **insesgado** de

$$r_i(\mathbf{x}) = \pi_i f(\mathbf{x}, \theta_i)$$

es decir,

$$\mathbb{E}_{\theta} \hat{r}_i(\mathbf{x}) = r_i(\mathbf{x}) \quad \text{para casi todo } \mathbf{x}$$

Entonces, si  $c_{ij} = 1$  para  $i \neq j$  se cumple que

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

## 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

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

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

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

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

## Estimación de los errores de clasificación

2) La tasa de error aparente es

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

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

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  $x \geq (x_1 + x_2)/2$  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 esta muestras que llamaremos  $\widehat{G}_0^{*,*}$  con regiones  $\widehat{G}_{i,0}^*$ . Sean

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

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

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

- 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

$$e_{i,boot} = \frac{m_i}{n_i} + \bar{d}_i \quad e_{boot} = \sum_{i=1}^k \pi_i e_{i,boot}$$

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





## 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  $k = 2$ , 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  $\pi_j = 1/k$  para todo  $j$ , la regla de clasificación Bayes es la obtenida antes con las coordenadas discriminantes.

- $\Sigma_B = \sum_{i=1}^k \pi_i (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})^T$ ,  $s = \text{rango}(\Sigma_B)$
- $\mathbf{z} = \mathbf{z} = \mathbf{A}^T \mathbf{x} = (\mathbf{z}_1^T, \mathbf{z}_2^T)^T$  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  $\mathbf{x}|G = i \sim N(\mu_i, \Sigma)$  vimos 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

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

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

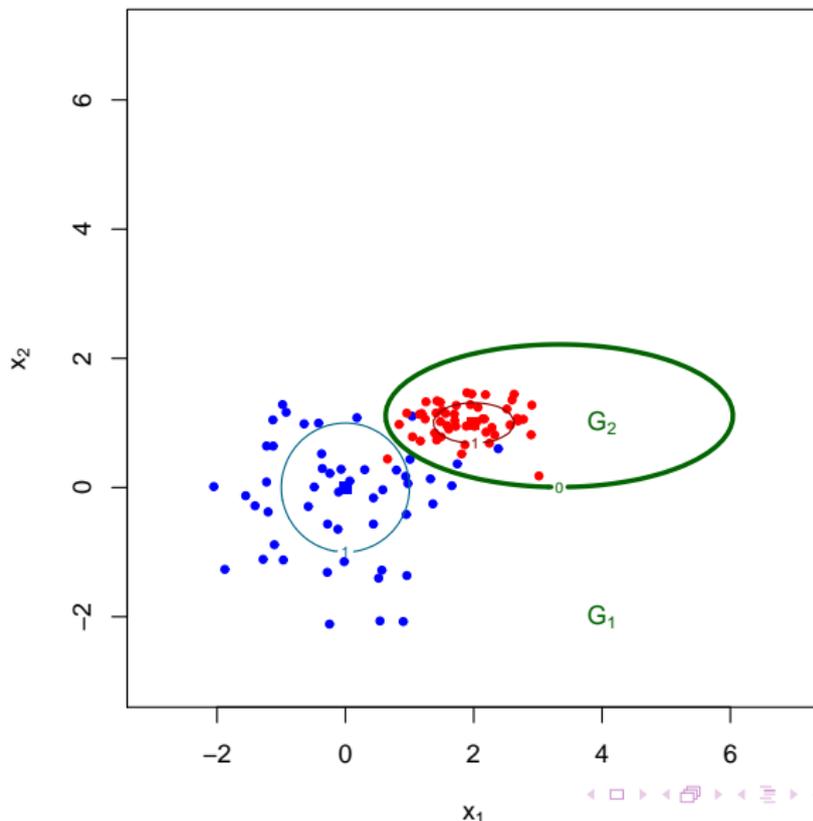
Para una probabilidad a priori  $\tau$  general tenemos que modificar  $\mathcal{G}_i$  por  $\mathcal{G}_{i,\tau}$ , o sea, asigno  $\mathbf{x}_0$  al grupo  $i$  si  $\mathbf{v}_0 = \mathbf{A}_1^T \mathbf{x}_0 \in \mathcal{G}_{i,\tau}$  donde

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



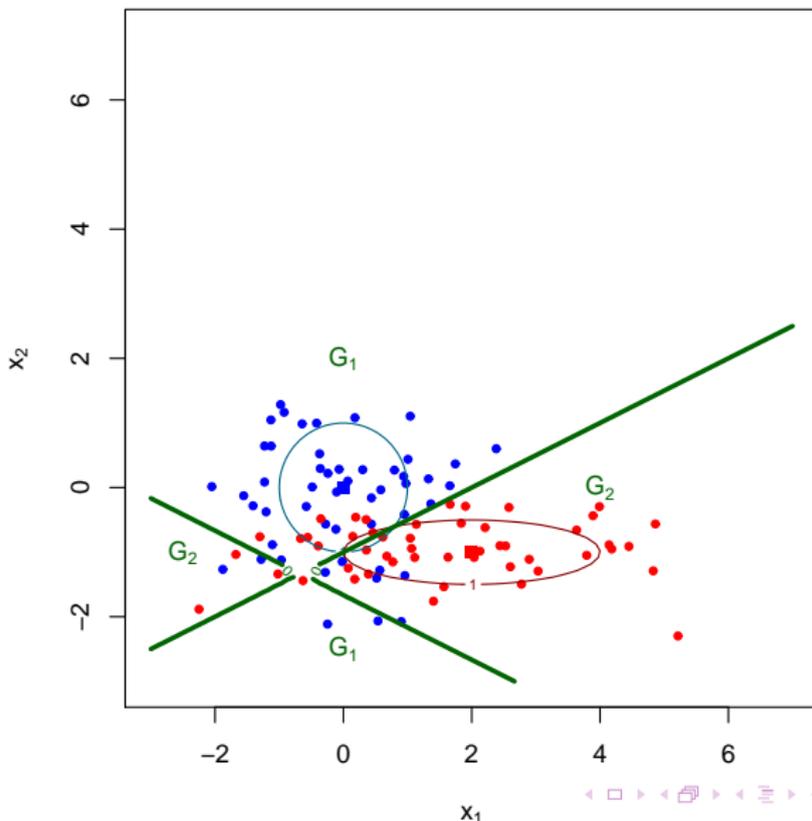


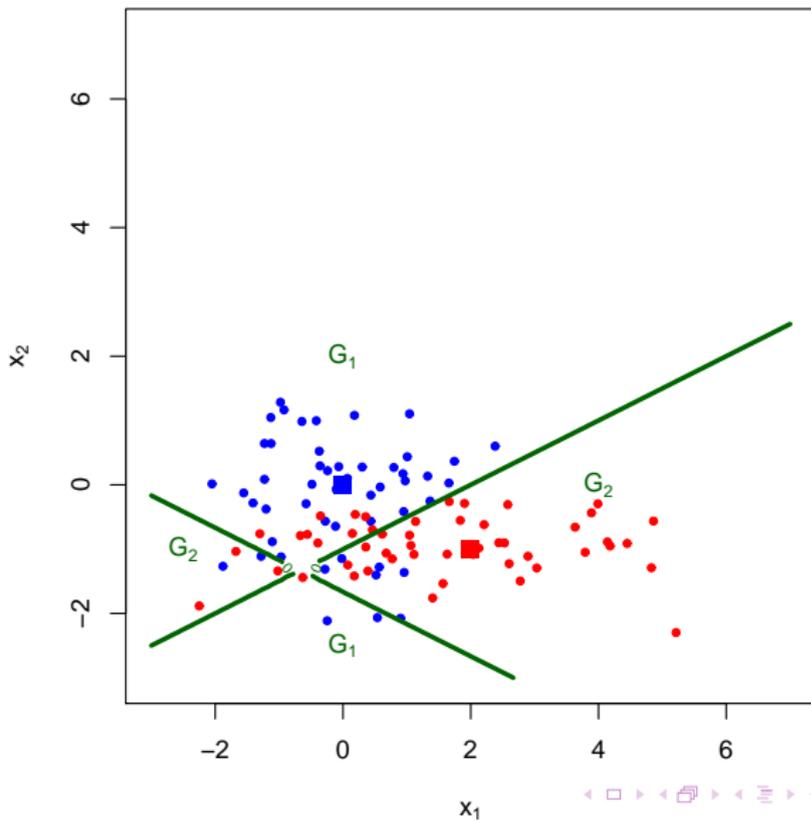


Ejemplo  $p = 2$ , Caso A

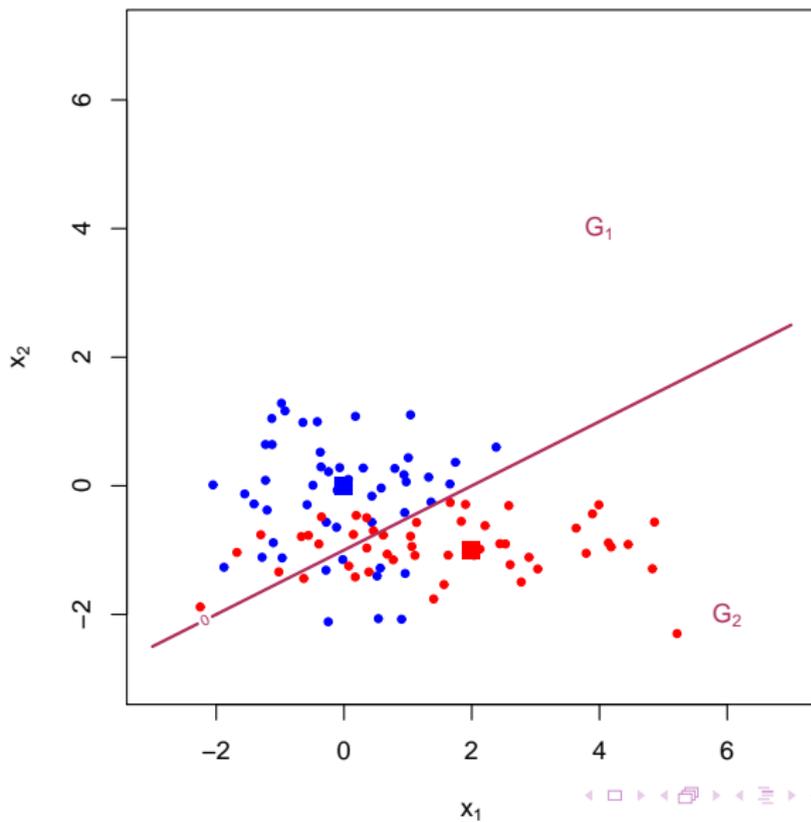


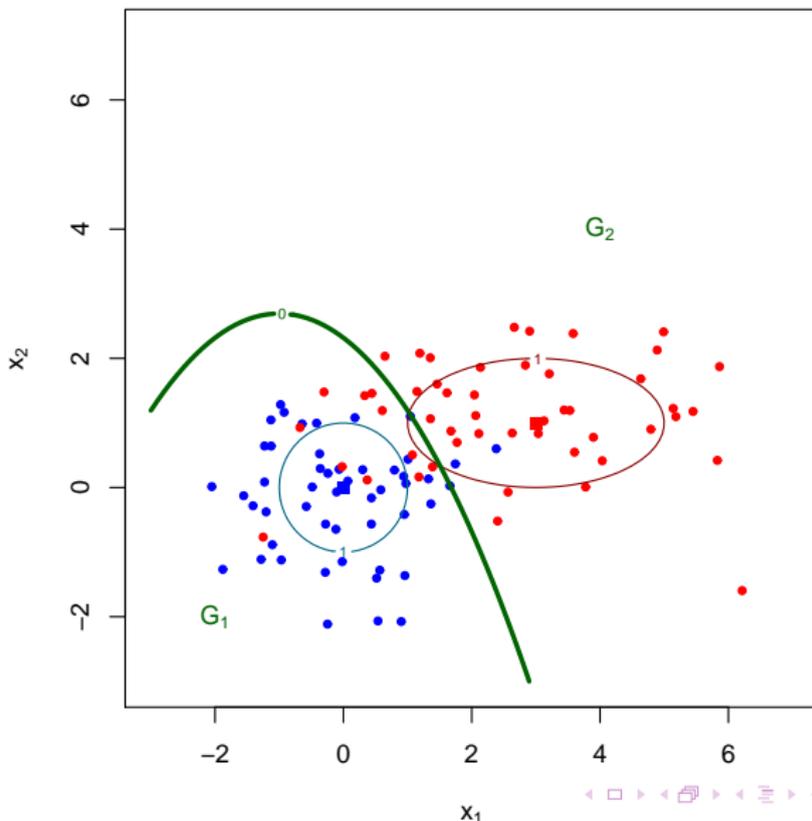


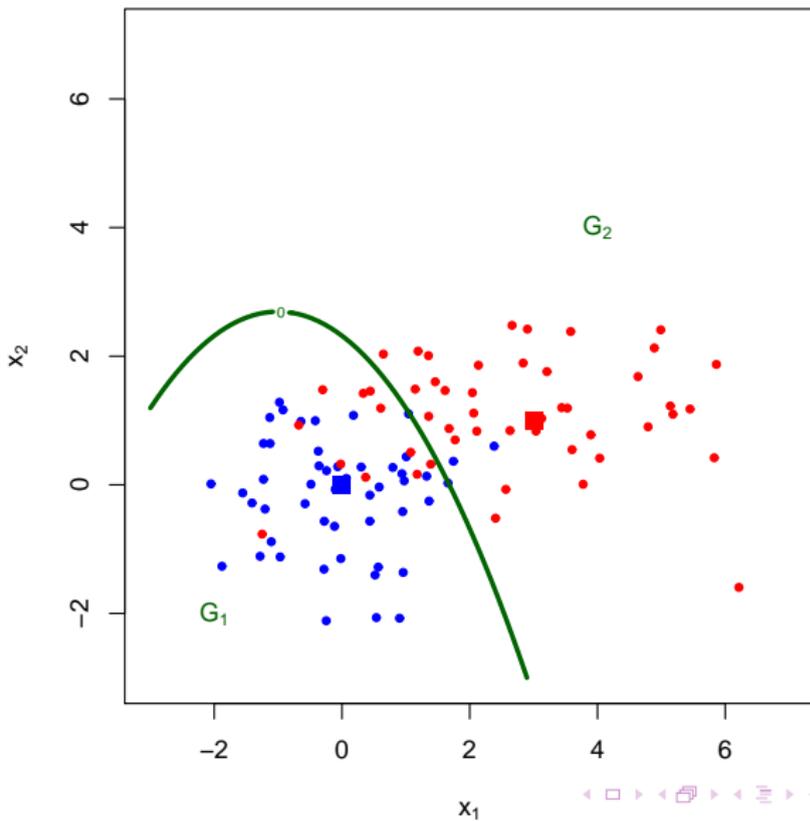
Ejemplo  $p = 2$ , Caso B

Ejemplo  $p = 2$ , Caso B

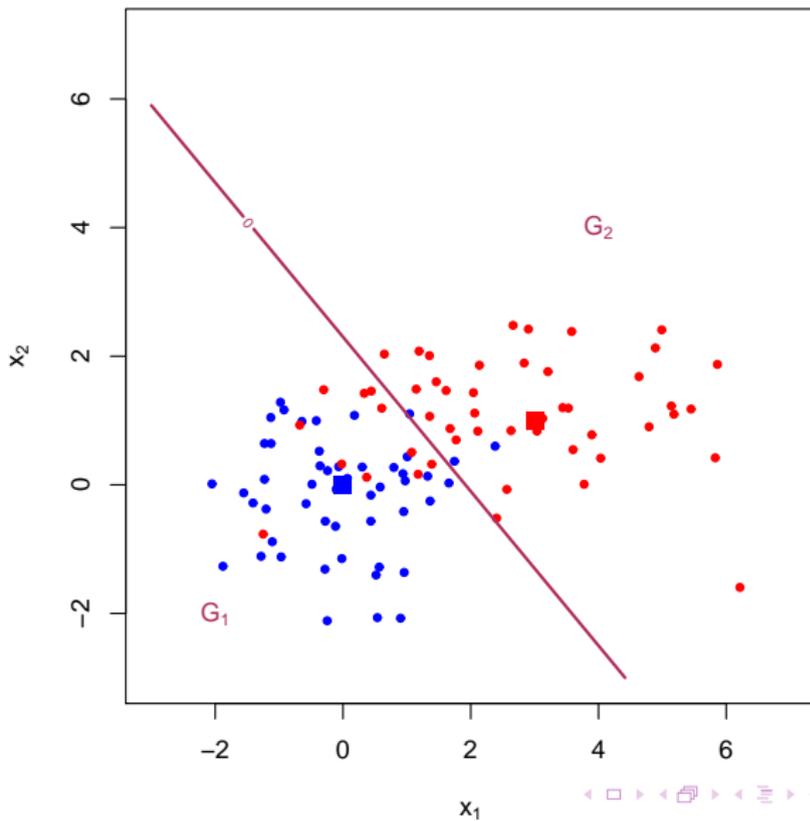
Ejemplo  $p = 2$ , Caso B: Regla Lineal como si  $\Sigma_1 = \Sigma_2$  lo cual es FALSO

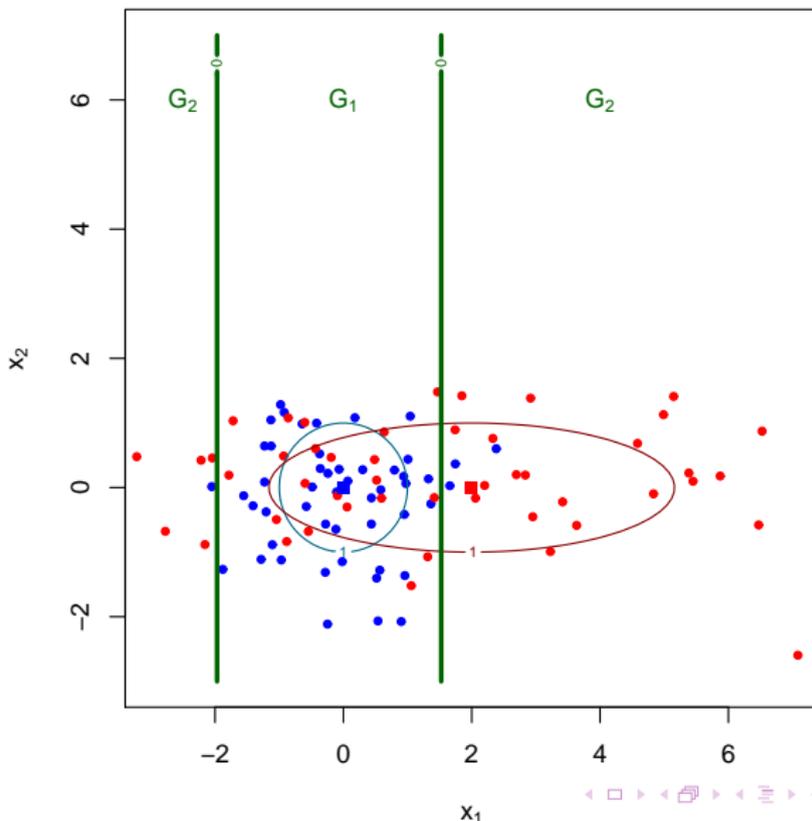


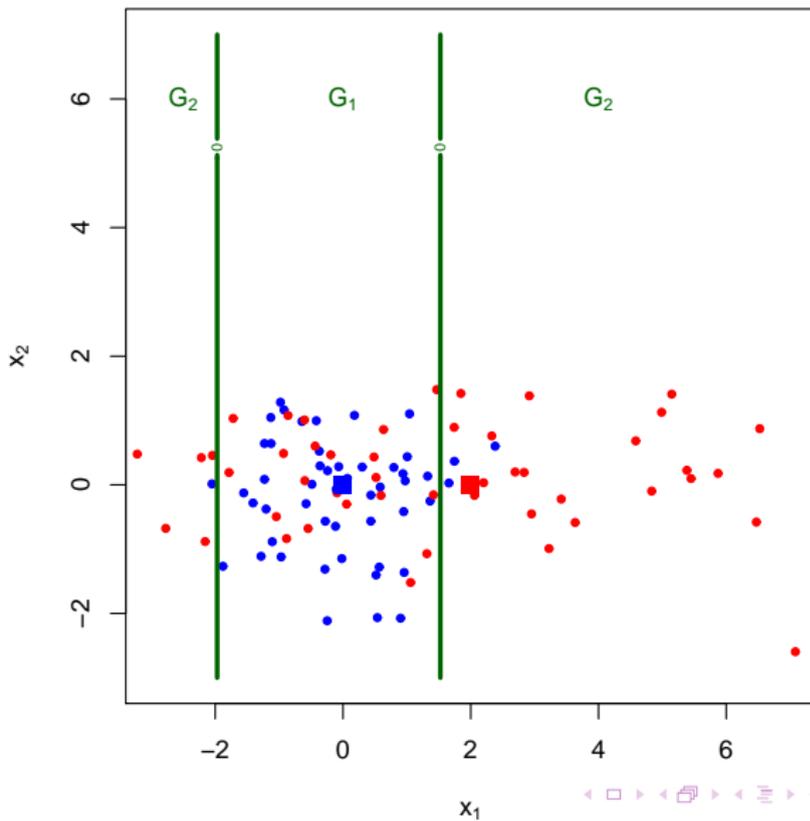
Ejemplo  $p = 2$ , Caso C,  $a_{22} = 0$ 

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

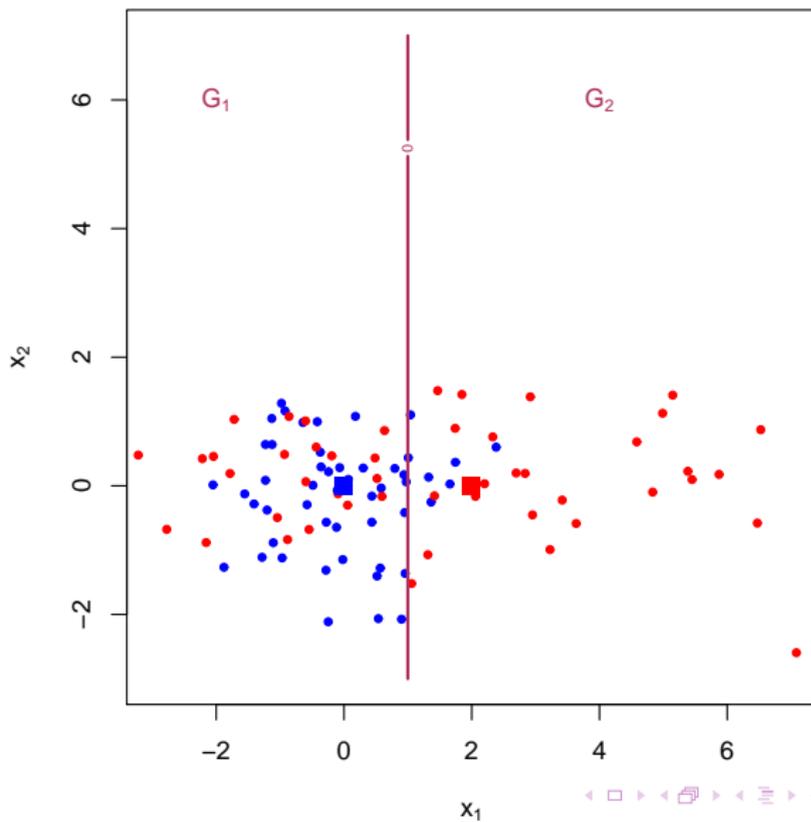
Ejemplo  $p = 2$ , Caso C: Regla Lineal como si  $\Sigma_1 = \Sigma_2$  lo cual es FALSO



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

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

Ejemplo  $p = 2$ , Caso D: Regla Lineal como si  $\Sigma_1 = \Sigma_2$  lo cual es FALSO







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

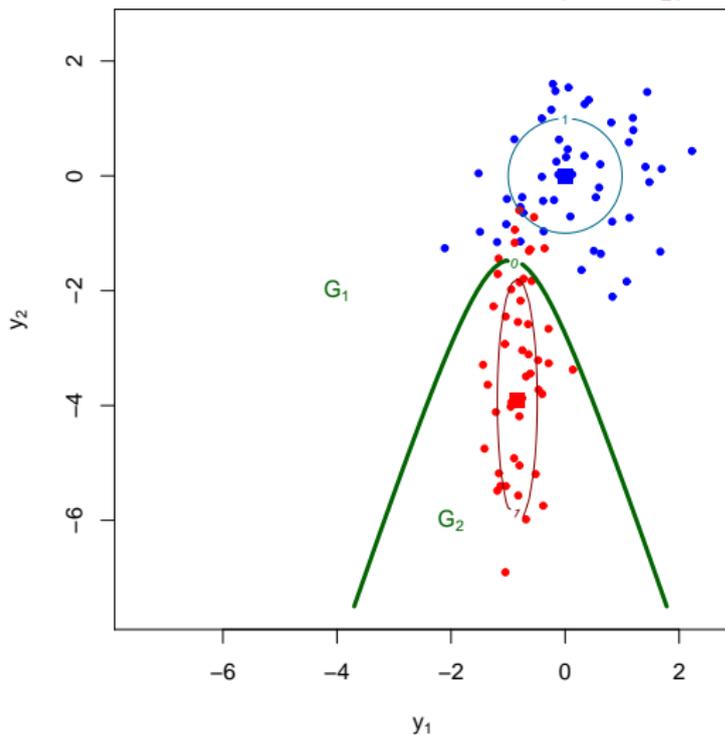
$$\mathbf{\Gamma} = \begin{pmatrix} -0.5590 & -1.1180 \\ -0.5590 & 1.1180 \end{pmatrix}$$

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

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

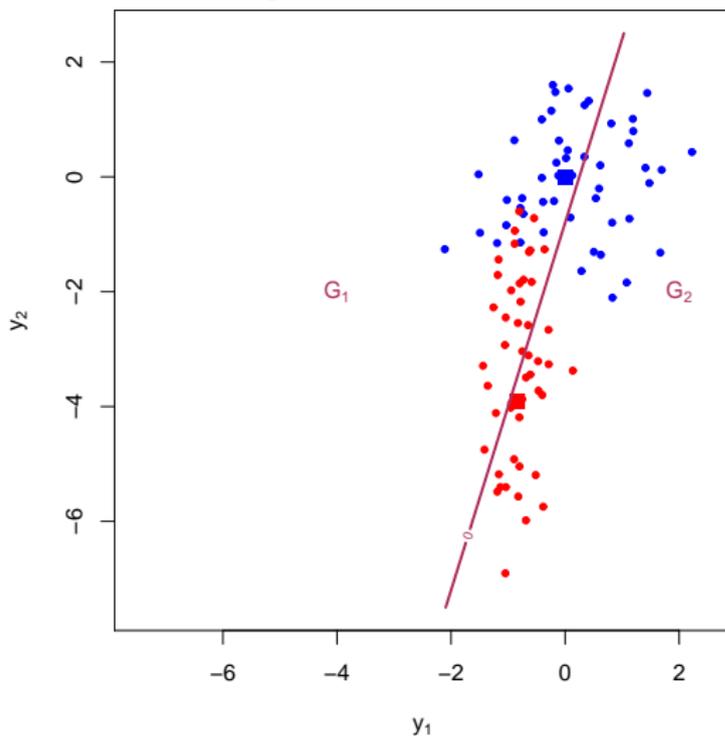
$$b_1 = 6.708 \quad b_2 = 0.8696 \quad c = 5.325$$

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





Datos transformados  $\mathbf{z} = \mathbf{\Gamma}^T(\mathbf{x} - \boldsymbol{\mu}_1)$ , Regla Lineal como si  $\boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$  lo cual es FALSO



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

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

$$\boldsymbol{\alpha}^T \mathbf{x} \sim N(\boldsymbol{\alpha}^T \boldsymbol{\mu}_1, \boldsymbol{\alpha}^T \boldsymbol{\Sigma} \boldsymbol{\alpha}) = N(\boldsymbol{\alpha}^T \boldsymbol{\mu}_1, \Delta_p^2)$$

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

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

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

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

$f_i \sim N(\boldsymbol{\mu}_i, \boldsymbol{\Sigma})$  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 $k = 2$ , $f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$

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

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

## 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  $k = 2$ ,  
 $f_j \sim N_p(\boldsymbol{\mu}_j, \boldsymbol{\Sigma})$

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

donde

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

Si  $\pi_1$  y  $\pi_2$  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  $\Delta_p^2$  y por lo tanto,  $\hat{e}_{act}$  subestima el error real  $e_{act}$ .

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

Esto se explica también pues recordemos que probamos el siguiente resultado

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

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







## 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  $H_0 : \boldsymbol{\Sigma}_1 = \boldsymbol{\Sigma}_2$ . 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  $\boldsymbol{\Sigma}_1$  y  $\boldsymbol{\Sigma}_2$ , LDF es preferible.
- Sin embargo, cuando  $p$  grande y  $\boldsymbol{\Sigma}_1$  y  $\boldsymbol{\Sigma}_2$  son muy distintas, la 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  $\Sigma_i$  si  $p$  es una fracción moderada de  $n_i$ .

## Modelos parsimoniosos: Ejemplo

Sean  $\mathbf{y}_{ij}$ ,  $1 \leq j \leq n_i$ ,  $i = 1, 2$  medidas correspondientes a hijos varones mellizos

- $n_1 = 49$  gemelos (monocigota)
- $n_2 = 40$  mellizos (heterocigota)

Se midieron

- altura,
- ancho de cadera,
- Circunferencia del pecho

del primer y segundo nacido, dando origen a un vector de dimensión 6.



## Modelos parsimoniosos: Ejemplo

Definamos

- $x_{ij,1} = y_{i,j,1} - y_{i,j,4}$   
= diferencia de altura entre los  $j$ -ésimos gemelos del grupo  $i$
- $x_{ij,2} = y_{i,j,2} - y_{i,j,5}$   
= diferencia de ancho de cadera entre los  $j$ -ésimos gemelos del grupo  $i$

Luego,

$$\mathbb{E}x_{ij} = \mathbf{0} \quad i = 1, 2, \quad 1 \leq j \leq n_i$$





## Modelos parsimoniosos

- $\mathbf{x}_i \sim N(\mathbf{0}_p, \rho_i \mathbf{\Sigma})$  con  $\rho_1 = 1, 1 = 1, 2$

La regla de discriminación en este caso es:

$$\mathcal{G}_{1,0} = \{\mathbf{x} : Q(\mathbf{x}) > 0\}$$

con  $Q(\mathbf{z}) = \mathbf{x}^T \mathbf{A} \mathbf{x} + c$

$$\mathbf{A} = \frac{1 - \rho}{2\rho} \mathbf{\Sigma}^{-1}$$

$$c = \log\left(\frac{\pi_1}{\pi_2}\right) + \frac{p}{2} \log(\rho)$$

## Modelos parsimoniosos

Sean  $\mathbf{x}_{ij} \sim N(\mathbf{0}_p, \rho_i \boldsymbol{\Sigma})$  con  $\rho_1 = 1$ ,  $1 \leq j \leq n_i$ ,  $i = 1, 2$ . Definamos

$$\mathbf{M}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} \mathbf{x}_{ij} \mathbf{x}_{ij}^T$$

Los EMV de  $\rho$  y  $\boldsymbol{\Sigma}$  resuelven

$$\hat{\boldsymbol{\Sigma}} = \frac{1}{n} \left( n_1 \mathbf{M}_1 + \frac{1}{\hat{\rho}} n_2 \mathbf{M}_2 \right)$$

$$\hat{\rho} = \frac{1}{p} \text{traza}(\hat{\boldsymbol{\Sigma}}^{-1} \mathbf{M}_2)$$

con lo cual

$$\hat{\mathbf{A}} = \frac{1 - \hat{\rho}}{2\hat{\rho}} \hat{\boldsymbol{\Sigma}}^{-1}$$

$$\hat{c} = \log \left( \frac{\pi_1}{\pi_2} \right) + \frac{p}{2} \log(\hat{\rho})$$

## Modelos parsimoniosos: Ejemplo

Sean  $\mathbf{y}_{ij}$ ,  $1 \leq j \leq n_i$ ,  $i = 1, 2$ .  $n_1 = 49$  son los gemelos y  $n_2 = 40$  los mellizos. Definimos:

- $x_{ij,1} = y_{i,j,1} - y_{i,j,4}$
- $x_{ij,2} = y_{i,j,2} - y_{i,j,5}$

$$\hat{\Sigma} = \begin{pmatrix} 3.8166 & 0.6462 \\ 0.6462 & 0.6473 \end{pmatrix}$$

$$\hat{\rho} = 7.070$$

# Modelos parsimoniosos: Ejemplo

mellizos                      gemelos

