

# MAT-269: Análisis Discriminante Lineal

**Felipe Osorio**

`fosorios.mat.utfsm.cl`

Departamento de Matemática, UTFSM



Suponga observaciones multivariadas provenientes de  $g$  clases (o grupos) **predefinidos**<sup>1</sup> teniendo características similares.

## *Ejemplos:*

Especies de plantas, Niveles de solvencia para clientes de un banco, presencia/ausencia de una condición médica, tipos de tumores, si un mensaje es SPAM o no, etc.

## Se desea:

- (a) **Discriminar**, esto es usar la información de aquellas **observaciones similares** para construir una **regla de clasificación** que permita separar tanto como sea posible las clases predefinidas.
- (b) **Clasificar**, dada las mediciones de una **nueva** observación **predecir** a que clase pertenece.

---

<sup>1</sup>Determinar los grupos será revisado más adelante, dentro de técnicas de agrupamiento (o clustering)



En esta clase consideraremos ( $g = 2$ ) clases o grupos y nuestro objetivo será contruir **un único** clasificador para diferenciar entre clases.

Suponga una población  $\mathcal{P}$  particionada en 2 grupos denotados por  $\Pi_1$  y  $\Pi_2$ . Además, cada elemento de  $\mathcal{P}$  es clasificado **sólo** en una clase.

Las mediciones de una muestra son usadas para asignar observaciones futuras a las clases designadas.

El vector aleatorio  $\mathbf{x} = (x_1, \dots, x_p)^\top$  representa las  $p$  mediciones de un ítem, las que son escogidas por su habilidad para distinguir entre las 2 clases.



Sea  $q_1$  y  $q_2$  las proporciones de individuos en  $\Pi_1$  y  $\Pi_2$ , respectivamente.

Podemos tener 2 tipos de errores en la clasificación:

- ▶ Clasificar un individuo en  $\Pi_1$ , cuando pertenece a  $\Pi_2$ .
- ▶ Clasificar un individuo en  $\Pi_2$ , cuando pertenece a  $\Pi_1$ .

Esto lleva a los siguientes costos de clasificación:

- ▶  $C(2|1)$  costo de clasificar un individuo en  $\Pi_2$  cuando realmente pertenece a  $\Pi_1$ .
- ▶  $C(1|2)$  costo de clasificar un individuo en  $\Pi_1$  cuando realmente pertenece a  $\Pi_2$ .

Sea,

- ▶  $g_1(\mathbf{x})$  la densidad de  $\mathbf{x}$  cuando un individuo pertenece a  $\Pi_1$ .
- ▶  $g_2(\mathbf{x})$  la densidad de  $\mathbf{x}$  cuando un individuo pertenece a  $\Pi_2$ .



Suponga las regiones  $\Omega_1$  y  $\Omega_2 \subset \mathbb{R}^p$  tal que  $\Omega_1 \cup \Omega_2 = \mathbb{R}^p$  y  $\Omega_1 \cap \Omega_2 = \emptyset$

De este modo la regla de decisión adopta la forma:

- ▶ Si  $x \in \Omega_1$ , clasificamos  $x$  en  $\Pi_1$ .
- ▶ Si  $x \in \Omega_2$ , clasificamos  $x$  en  $\Pi_2$ .

## Objetivo:

Definir regiones  $\Omega_1$  y  $\Omega_2$  que minimicen los costos esperados de clasificación errónea.

Sea,

- ▶  $P(2|1)$  probabilidad de clasificar erróneamente un individuo de  $\Pi_1$  como perteneciente a  $\Pi_2$ .
- ▶  $P(1|2)$  probabilidad de clasificar erróneamente un individuo de  $\Pi_2$  como perteneciente a  $\Pi_1$ .



Tenemos,

$$P(2|1) = P\{x \in \Omega_2 \text{ cuando } x \sim g_1(x)\} = \int_{\Omega_2} g_1(x) dx,$$

y análogamente

$$P(1|2) = P\{x \in \Omega_1 \text{ cuando } x \sim g_2(x)\} = \int_{\Omega_1} g_2(x) dx.$$

Suponga que la probabilidad de obtener una observación desde  $\Pi_1$  es  $q_1$  y análogamente la probabilidad de obtener una observación desde  $\Pi_2$  es  $q_2$ . Además,  $q_1 + q_2 = 1$ .

De este modo, la probabilidad de que un individuo proveniente de  $\Pi_1$  ( $\Pi_2$ ) sea clasificado erróneamente es  $q_1 P(2|1)$  ( $q_2 P(1|2)$ ).

Esto lleva a la siguiente tabla:

Costo	$C(2 1)$	$C(1 2)$
Probabilidad	$q_1 P(2 1)$	$q_2 P(1 2)$



Luego

$$\begin{aligned} E(\text{Clasificación}) &= C(2|1)q_1 P(2|1) + C(1|2)q_2 P(1|2) \\ &= C_1 \int_{\Omega_2} g_1(\mathbf{x})d\mathbf{x} + C_2 \int_{\Omega_1} g_2(\mathbf{x})d\mathbf{x}, \end{aligned}$$

donde  $C_1 = q_1 C(2|1)$  y  $C_2 = q_2 C(1|2)$ . Luego,

$$\begin{aligned} E(\text{Clasificación}) &= C_1 \left[ 1 - \int_{\Omega_1} g_1(\mathbf{x})d\mathbf{x} \right] + C_2 \int_{\Omega_1} g_2(\mathbf{x})d\mathbf{x} \\ &= C_1 + \int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\}d\mathbf{x}. \end{aligned}$$

Se desea determinar  $\Omega_1$  minimizando la integral:

$$\int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\}d\mathbf{x}$$



Por notar que  $\Omega_1 = \{\mathbf{x} : C_1 g_1(\mathbf{x}) \geq C_2 g_2(\mathbf{x})\}$  minimiza la integral

$$\int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x}$$

Recuerde que:

$$I_{\Omega_1}(\mathbf{x}) = \begin{cases} 1, & \mathbf{x} \in \Omega_1, \\ 0, & \text{en otro caso.} \end{cases}$$

Lleva a escribir

$$\int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x} = \int I_{\Omega_1}(\mathbf{x}) \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x}.$$





Sea  $\Omega_1^*$  otra región. Se desea mostrar que

$$\int_{\Omega_1^*} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x} - \int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x} \geq 0$$

Lo que es equivalente a mostrar que

$$\int \{I_{\Omega_1^*}(\mathbf{x}) - I_{\Omega_1}(\mathbf{x})\} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x} \geq 0$$

Es decir, basta probar que:

$$\{I_{\Omega_1^*}(\mathbf{x}) - I_{\Omega_1}(\mathbf{x})\} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} \geq 0$$



Supongamos que  $\{I_{\Omega_1^*}(\mathbf{x}) - I_{\Omega_1}(\mathbf{x})\} > 0$ . Entonces,  $I_{\Omega_1}(\mathbf{x}) = 0$ . Esto implica que  $\mathbf{x} \notin \Omega_1$ .

Si

$$\{I_{\Omega_1^*}(\mathbf{x}) - I_{\Omega_1}(\mathbf{x})\} < 0,$$

entonces  $I_{\Omega_1}(\mathbf{x}) = 1 \Rightarrow \mathbf{x} \in \Omega_1$ . Esto implica que

$$\{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} \leq 0.$$

Por simplicidad asuma que  $C(1|2) = C(2|1)$  y que  $q_2 = q_1$ . Entonces la condición anterior resulta,

$$g_2(\mathbf{x}) \leq g_1(\mathbf{x}) \quad \text{o bien,} \quad \frac{g_1(\mathbf{x})}{g_2(\mathbf{x})} \geq 1.$$

Así, finalmente podemos escribir la región que minimiza la integral deseada como:

$$\Omega_1 = \left\{ \mathbf{x} : \frac{g_1(\mathbf{x})}{g_2(\mathbf{x})} \geq 1 \right\}.$$



### Ejemplo:

Considere  $g_1(\mathbf{x}) \stackrel{d}{=} N_p(\boldsymbol{\mu}_1, \boldsymbol{\Sigma})$  y  $g_2(\mathbf{x}) \stackrel{d}{=} N_p(\boldsymbol{\mu}_2, \boldsymbol{\Sigma})$ .

Entonces

$$\frac{g_1(\mathbf{x})}{g_2(\mathbf{x})} = \frac{\exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_1)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_1))}{\exp(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}_2))}$$

Luego, la región  $\Omega_1$  está determinada por la condición

$$\log\left(\frac{g_1(\mathbf{x})}{g_2(\mathbf{x})}\right) \geq \log K, \quad K = \frac{C(1|2)}{C(2|1)} \frac{q_2}{q_1}.$$

Ahora

$$\log\left(\frac{g_1(\mathbf{x})}{g_2(\mathbf{x})}\right) = \mathbf{x}^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Por lo tanto,

$$\Omega_1 = \left\{ \mathbf{x} : \mathbf{x}^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \geq \log K \right\}.$$



## Análisis discriminante

Si  $C(1|2) = C(2|1)$  y  $q_1 = q_2$ , obtenemos la condición:

$$\Omega_1 = \left\{ \mathbf{x} : \mathbf{x}^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \geq \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \right\}.$$

La función  $L(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$  se denomina **función discriminante lineal de Fisher**.

Para evaluar  $P(2|1)$  y  $P(1|2)$  definamos la función:

$$U(\mathbf{x}) = \mathbf{x}^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2).$$

Entonces, de acuerdo al desarrollo anterior, la región  $\Omega_1$  queda determinada por la condición  $U(\mathbf{x}) \geq 0$ .

Asumamos que  $U(\mathbf{x})$  tiene una distribución normal univariada. Entonces, si  $\mathbf{x}$  pertenece a  $\Pi_1$ ,

$$\begin{aligned} E\{U(\mathbf{x})\} &= \boldsymbol{\mu}_1^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) - \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= \frac{1}{2}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \Sigma^{-1}(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = \frac{1}{2}\Delta^2. \end{aligned}$$



## Análisis discriminante

Note que

$$\Delta^2 = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

es el cuadrado de la **distancia de Mahalanobis** entre las dos medias poblacionales.

Ahora note que

$$\begin{aligned}\text{var}\{U(\boldsymbol{x})\} &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} \boldsymbol{\Sigma} \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) \\ &= (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2) = \Delta^2\end{aligned}$$

Luego, si  $\boldsymbol{x}$  proviene de  $\Pi_1$ , entonces

$$U(\boldsymbol{x}) \sim N(\tfrac{1}{2}\Delta^2, \Delta^2).$$

Similarmente se puede probar que si  $\boldsymbol{x}$  proviene de  $\Pi_2$ , entonces

$$U(\boldsymbol{x}) \sim N(-\tfrac{1}{2}\Delta^2, \Delta^2).$$

Finalmente para evaluar  $P(1|2)$  calculamos

$$P(1|2) = P\{U(\boldsymbol{x}) \geq 0 \text{ cuando } U(\boldsymbol{x}) \sim N(-\tfrac{1}{2}\Delta^2, \Delta^2)\}.$$



- ▶ Si  $q_1$  y  $q_2$  no están disponibles, es posible derivar una regla de clasificación poniendo condiciones sobre  $C(1|2)$ ,  $C(2|1)$ ,  $P(1|2)$ , y  $P(2|1)$ .
- ▶ Por ejemplo, si ambas poblaciones son normales, digamos,  $N(\mu_1, \Sigma)$  y  $N(\mu_2, \Sigma)$  y la región  $\Omega_1$  es determinada por la condición  $U(x) \geq c$ .
- ▶ Podemos calcular

$$\begin{aligned} P(2|1) &= P\{U(x) < c \text{ cuando } x \sim N(\mu_1, \Sigma)\} \\ &= P\{U(x) < c \text{ cuando } U(x) \sim N(\tfrac{1}{2}\Delta^2, \Delta^2)\} \\ &= \int_{-\infty}^{\frac{c - \Delta^2/2}{\Delta}} \frac{2}{\sqrt{2\pi}} \exp(-\tfrac{1}{2}y^2) dy. \end{aligned}$$

Luego podemos determinar una constante  $c$  tal que, por ejemplo,  $P(2|1) = 0.05$ .

- ▶ Otro tipo de restricciones que pueden considerarse son, por ejemplo,

$$C(1|2) P(1|2) = C(2|1) P(2|1).$$



### Resultado:

En el contexto del problema de minimización de la integral

$$\int_{\Omega_1} \{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} d\mathbf{x},$$

supongamos que para otra región  $\Omega_1^*$  que satisface

$$(I_{\Omega_1^*}(\mathbf{x}) - I_{\Omega_1}(\mathbf{x}))\{C_2 g_2(\mathbf{x}) - C_1 g_1(\mathbf{x})\} \geq 0$$

se tiene que  $\Omega_1^*$  también minimiza el costo esperado. Entonces

$$\Omega_1^* = \Omega_1.$$

### Demostración:

La demostración consiste en probar que  $I_{\Omega_1^*}(\mathbf{x}) = I_{\Omega_1}(\mathbf{x})$ .

