

# Análisis Factorial Discriminante y Análisis Discriminante

Grado en Matemáticas — Análisis de Datos

Norberto Corral Blanco

Beatriz Sinova Fernández

Universidad de Oviedo

Universidad de Oviedo

## Índice

<b>1. Preliminares: matrices de datos y centrado</b>	<b>2</b>
1.1. La matriz de datos . . . . .	2
1.2. Matriz de centrado . . . . .	2
<b>2. Sumas de cuadrados y descomposición de la variabilidad</b>	<b>3</b>
2.1. Suma de cuadrados total . . . . .	3
2.2. Estructura por grupos . . . . .	3
2.3. Descomposición $\mathbf{T} = \mathbf{W} + \mathbf{B}$ . . . . .	3
2.3.1. Variabilidad entre grupos: $\mathbf{B}$ . . . . .	3
2.3.2. Variabilidad intragrupos: $\mathbf{W}$ . . . . .	4
<b>3. Análisis Factorial Discriminante: factores canónicos</b>	<b>4</b>
3.1. Motivación . . . . .	4
3.2. Planteamiento del problema . . . . .	5
3.3. Resolución: reducción a un problema de valores propios . . . . .	5
3.3.1. Conexión con los valores propios de $\mathbf{W}^{-1}\mathbf{B}$ . . . . .	5
3.4. Número de factores discriminantes canónicos . . . . .	6
3.5. Los $s$ factores discriminantes son no correlacionados . . . . .	6
3.6. Interpretación de los factores . . . . .	6
<b>4. Análisis Discriminante: clasificación</b>	<b>7</b>
4.1. El problema de clasificación . . . . .	7
4.2. El método de máxima verosimilitud . . . . .	7
4.2.1. Caso univariante con igual varianza . . . . .	7
4.2.2. Caso univariante con desiguales varianzas . . . . .	8
4.3. Caso multivariante: discriminante lineal (LDA) . . . . .	8
4.3.1. Dos subpoblaciones . . . . .	8
4.3.2. $g$ subpoblaciones con igual $\Sigma$ . . . . .	8
4.4. Discriminante cuadrático (QDA) . . . . .	9
4.5. El método bayesiano . . . . .	9
4.5.1. LDA bayesiano . . . . .	9
4.5.2. QDA bayesiano . . . . .	9

5. Estimación de parámetros desconocidos	10
6. Resumen y relación entre métodos	10

# 1. Preliminares: matrices de datos y centrado

## 1.1. La matriz de datos

Consideramos una **matriz de datos**  $\mathbf{X}$  de dimensión  $n \times p$ , donde  $n$  es el número de observaciones y  $p$  el número de variables:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \cdots & x_{1p} \\ x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n1} & x_{n2} & \cdots & x_{np} \end{pmatrix}.$$

Cada fila  $\mathbf{x}_i^t = (x_{i1}, \dots, x_{ip})$  representa la observación  $i$ -ésima; cada columna  $\mathbf{x}_{(j)}$  recoge los  $n$  valores de la variable  $j$ .

## 1.2. Matriz de centrado

**Definición 1.1.** La **matriz de centrado** es

$$\mathbf{H} = \mathbf{I}_n - \frac{1}{n} \mathbf{1}\mathbf{1}^t = \mathbf{I}_n - \mathbf{J}, \quad \mathbf{J} = \frac{1}{n} \mathbf{1}\mathbf{1}^t.$$

La acción de  $\mathbf{H}$  sobre un vector  $\mathbf{x}$  es restar la media:

$$\mathbf{H}\mathbf{x} = \mathbf{x} - \mathbf{1}\bar{x} = (x_1 - \bar{x}, \dots, x_n - \bar{x})^t.$$

La **matriz de datos centrada** respecto a todas las variables es

$$\mathbf{X}_c = \mathbf{H}\mathbf{X},$$

cuya entrada  $(i, j)$  vale  $x_{ij} - \bar{x}_j$ .

**Proposición 1.2** (Propiedades de  $\mathbf{H}$ ). (I)  $\mathbf{H}$  es *simétrica*:  $\mathbf{H}^t = \mathbf{H}$ .

(II)  $\mathbf{H}$  es *idempotente*:  $\mathbf{H}^2 = \mathbf{H}$ .

(III) Los valores propios de  $\mathbf{H}$  son 0 (con multiplicidad 1) y 1 (con multiplicidad  $n - 1$ ).

(IV)  $\text{rango}(\mathbf{H}) = \text{tr}(\mathbf{H}) = n - 1$ .

*Demostración.* (i)  $\mathbf{H}^t = \mathbf{I}_n^t - \frac{1}{n}(\mathbf{1}\mathbf{1}^t)^t = \mathbf{I}_n - \frac{1}{n}\mathbf{1}\mathbf{1}^t = \mathbf{H}$ .

(ii)  $\mathbf{J}^2 = \frac{1}{n^2}\mathbf{1}\mathbf{1}^t\mathbf{1}\mathbf{1}^t = \frac{n}{n^2}\mathbf{1}\mathbf{1}^t = \mathbf{J}$ , por lo que

$$\mathbf{H}^2 = (\mathbf{I} - \mathbf{J})^2 = \mathbf{I} - 2\mathbf{J} + \mathbf{J}^2 = \mathbf{I} - 2\mathbf{J} + \mathbf{J} = \mathbf{I} - \mathbf{J} = \mathbf{H}.$$

(iii) Si  $\mathbf{H}\mathbf{u} = \lambda\mathbf{u}$ , aplicando  $\mathbf{H}$ :  $\mathbf{H}^2\mathbf{u} = \lambda^2\mathbf{u}$ ; pero  $\mathbf{H}^2 = \mathbf{H}$ , así que  $\lambda\mathbf{u} = \lambda^2\mathbf{u}$ , es decir  $\lambda(\lambda - 1) = 0$ . El único vector en el núcleo de  $\mathbf{H}$  es proporcional a  $\mathbf{1}$  (pues  $\mathbf{H}\mathbf{1} = \mathbf{1} - \mathbf{1} = \mathbf{0}$ ), de modo que  $\lambda = 0$  tiene multiplicidad 1 y  $\lambda = 1$  tiene multiplicidad  $n - 1$ .

(iv) Como  $\mathbf{H}$  es simétrica e idempotente, admite diagonalización ortogonal  $\mathbf{H} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^t$  con  $\mathbf{U}$  ortogonal. Por tanto  $\text{rango}(\mathbf{H}) = \text{rango}(\mathbf{\Lambda}) = \text{tr}(\mathbf{\Lambda}) = \text{tr}(\mathbf{U}\mathbf{\Lambda}\mathbf{U}^t) = \text{tr}(\mathbf{H}) = n - 1$ .  $\square$

## 2. Sumas de cuadrados y descomposición de la variabilidad

### 2.1. Suma de cuadrados total

**Definición 2.1.** La suma de cuadrados total es la matriz  $p \times p$

$$\mathbf{T} = \mathbf{X}^t \mathbf{H} \mathbf{X}.$$

La matriz de varianzas-covarianzas muestral verifica  $\mathbf{S} = \frac{1}{n} \mathbf{T}$ . Cuando  $n > p$ , el estimador insesgado de la matriz de covarianzas poblacional es  $\hat{\mathbf{S}} = \frac{1}{n-p} \mathbf{T}$ .

### 2.2. Estructura por grupos

Supongamos que las  $n$  observaciones provienen de  $g$  subpoblaciones: el grupo  $k$  aporta  $n_k$  observaciones ( $\sum_{k=1}^g n_k = n$ ). La matriz de datos se particiona como  $\mathbf{X}^t = (\mathbf{X}_1^t \mid \cdots \mid \mathbf{X}_g^t)$ , donde  $\mathbf{X}_k$  tiene dimensión  $n_k \times p$ .

**Definición 2.2.** Definimos la **matriz diagonal por bloques**

$$\mathbf{D} = \text{diag}(\mathbf{J}_1, \dots, \mathbf{J}_g), \quad \mathbf{J}_k = \frac{1}{n_k} \mathbf{1}_{n_k} \mathbf{1}_{n_k}^t.$$

**Proposición 2.3** (Propiedades de  $\mathbf{D}$ ).  $\mathbf{D}$  es simétrica e idempotente con  $\text{rango}(\mathbf{D}) = \text{tr}(\mathbf{D}) = g$ .

*Demostración.* Simetría: cada  $\mathbf{J}_k$  es simétrica, luego  $\mathbf{D}^t = \mathbf{D}$ . Idempotencia:  $\mathbf{D}^2 = \text{diag}(\mathbf{J}_1^2, \dots, \mathbf{J}_g^2) = \text{diag}(\mathbf{J}_1, \dots, \mathbf{J}_g) = \mathbf{D}$  porque  $\mathbf{J}_k^2 = \mathbf{J}_k$ . Rango:  $\text{tr}(\mathbf{J}_k) = n_k \cdot \frac{1}{n_k} = 1$ , así que  $\text{rango}(\mathbf{D}) = \text{tr}(\mathbf{D}) = \sum_{k=1}^g 1 = g$ .  $\square$

### 2.3. Descomposición $\mathbf{T} = \mathbf{W} + \mathbf{B}$

**Teorema 2.4** (Descomposición de la variabilidad total).

$$\mathbf{T} = \mathbf{X}^t \mathbf{H} \mathbf{X} = \underbrace{\mathbf{X}^t (\mathbf{I} - \mathbf{D}) \mathbf{X}}_{\mathbf{W}} + \underbrace{\mathbf{X}^t (\mathbf{D} - \mathbf{J}) \mathbf{X}}_{\mathbf{B}}.$$

*Demostración.* Basta escribir  $\mathbf{H} = \mathbf{I} - \mathbf{J} = (\mathbf{I} - \mathbf{D}) + (\mathbf{D} - \mathbf{J})$ .  $\square$

La **variabilidad intragrupos** ( $\mathbf{W}$ ) y la **variabilidad entre grupos** ( $\mathbf{B}$ ) admiten las siguientes expresiones explícitas.

#### 2.3.1. Variabilidad entre grupos: $\mathbf{B}$

Observemos que  $\mathbf{D}\mathbf{J} = \mathbf{J}$  (cada bloque  $\mathbf{J}_k$  aplicado a  $\mathbf{J}$  da  $\mathbf{J}$  por ser todas las columnas de  $\mathbf{J}$  iguales a  $\frac{1}{n}$ ). Por simetría,  $\mathbf{J}\mathbf{D} = \mathbf{J}$  también.

**Proposición 2.5** (Propiedades de  $\mathbf{D} - \mathbf{J}$ ).  $\mathbf{D} - \mathbf{J}$  es simétrica e idempotente con  $\text{rango}(\mathbf{D} - \mathbf{J}) = g - 1$ .

*Demostración. Simetría:* inmediata por serlo  $\mathbf{D}$  y  $\mathbf{J}$ . *Idempotencia:*

$$(\mathbf{D} - \mathbf{J})^2 = \mathbf{D}^2 - \mathbf{D}\mathbf{J} - \mathbf{J}\mathbf{D} + \mathbf{J}^2 = \mathbf{D} - \mathbf{J} - \mathbf{J} + \mathbf{J} = \mathbf{D} - \mathbf{J}.$$

*Rango:*  $\text{rango}(\mathbf{D} - \mathbf{J}) = \text{tr}(\mathbf{D} - \mathbf{J}) = \text{tr}(\mathbf{D}) - \text{tr}(\mathbf{J}) = g - 1.$   $\square$

Dado que  $(\mathbf{D} - \mathbf{J})$  es simétrica e idempotente, usando que  $\mathbf{D}\mathbf{X} = (\mathbf{J}_1\mathbf{X}_1^t \mid \cdots \mid \mathbf{J}_g\mathbf{X}_g^t)^t$  con  $\mathbf{J}_k\mathbf{X}_k$  igual a la matriz cuyas filas son  $\bar{\mathbf{x}}_k^t$ , se obtiene:

$$\mathbf{B} = \mathbf{X}^t(\mathbf{D} - \mathbf{J})\mathbf{X} = [(\mathbf{D} - \mathbf{J})\mathbf{X}]^t [(\mathbf{D} - \mathbf{J})\mathbf{X}] = \sum_{k=1}^g n_k (\bar{\mathbf{x}}_k - \bar{\mathbf{x}})(\bar{\mathbf{x}}_k - \bar{\mathbf{x}})^t,$$

donde  $\bar{\mathbf{x}}_k$  es el vector de medias del grupo  $k$  y  $\bar{\mathbf{x}}$  la media global.

### 2.3.2. Variabilidad intragrupos: $\mathbf{W}$

**Proposición 2.6** (Propiedades de  $\mathbf{I} - \mathbf{D}$ ).  $\mathbf{I} - \mathbf{D}$  es simétrica e idempotente con  $\text{rango}(\mathbf{I} - \mathbf{D}) = n - g.$

*Demostración.* Análogo al caso de  $\mathbf{H}$ . Idempotencia:  $(\mathbf{I} - \mathbf{D})^2 = \mathbf{I} - 2\mathbf{D} + \mathbf{D}^2 = \mathbf{I} - \mathbf{D}$ . Rango:  $\text{tr}(\mathbf{I} - \mathbf{D}) = n - g.$   $\square$

La variabilidad intragrupos satisface:

$$\mathbf{W} = \mathbf{X}^t(\mathbf{I} - \mathbf{D})\mathbf{X} = \sum_{k=1}^g \sum_{i=1}^{n_k} (\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)(\mathbf{x}_{ki} - \bar{\mathbf{x}}_k)^t = \sum_{k=1}^g n_k \mathbf{S}_k,$$

donde  $\mathbf{x}_{ki}$  es la  $i$ -ésima observación del grupo  $k$  y  $\mathbf{S}_k$  es la matriz de covarianzas muestral del grupo  $k$ .

## 3. Análisis Factorial Discriminante: factores canónicos

### 3.1. Motivación

Clásicamente, el análisis discriminante distingue dos vertientes:

- **Aspecto descriptivo (geométrico):** buscar combinaciones lineales de las variables originales que separen lo mejor posible los  $g$  grupos.
- **Aspecto decisonal (probabilístico):** asignar un nuevo individuo a uno de los grupos a partir de sus valores observados.

A menudo, cada variable marginal  $x_j$  se comporta de modo similar en todos los grupos, pero las diferencias resultan muy claras al considerar el comportamiento conjunto. Los **factores discriminantes canónicos** responden al aspecto descriptivo.

### 3.2. Planteamiento del problema

Dada la combinación lineal centrada  $\mathbf{y} = \mathbf{H}\mathbf{X}\mathbf{a}$ , su variabilidad total se descompone como:

$$n \operatorname{Var}(\mathbf{y}) = \mathbf{y}^t \mathbf{y} = \mathbf{a}^t \mathbf{X}^t \mathbf{H} \mathbf{X} \mathbf{a} = \mathbf{a}^t \mathbf{T} \mathbf{a} = \underbrace{\mathbf{a}^t \mathbf{W} \mathbf{a}}_{\text{intragrupos}} + \underbrace{\mathbf{a}^t \mathbf{B} \mathbf{a}}_{\text{entre grupos}}.$$

**Objetivo:** Encontrar  $\mathbf{a} \in \mathbb{R}^p$  que maximice la proporción de variabilidad explicada entre grupos respecto a la variabilidad intragrupos:

$$\max_{\mathbf{a}} \frac{\mathbf{a}^t \mathbf{B} \mathbf{a}}{\mathbf{a}^t \mathbf{W} \mathbf{a}}.$$

Maximizar únicamente  $\mathbf{a}^t \mathbf{B} \mathbf{a}$  no tiene sentido: una dirección con mayor variabilidad entre grupos puede tener aún mayor dispersión intragrupos, lo que empeoraría la separación real.

### 3.3. Resolución: reducción a un problema de valores propios

Suponemos que  $\mathbf{W}$  es **definida positiva** (lo cual requiere  $n - g \geq p$  y que ninguna variable sea combinación lineal de las restantes). En ese caso,  $\mathbf{W} = \mathbf{U}\mathbf{\Lambda}\mathbf{U}^t$  con  $\mathbf{\Lambda}$  diagonal de entradas positivas, y podemos definir

$$\mathbf{W}^{1/2} = \mathbf{U}\mathbf{\Lambda}^{1/2}\mathbf{U}^t, \quad \mathbf{W}^{-1/2} = \mathbf{U}\mathbf{\Lambda}^{-1/2}\mathbf{U}^t.$$

Ambas son **simétricas** e inversas la una de la otra. El cambio de variable  $\mathbf{b} = \mathbf{W}^{1/2}\mathbf{a}$  transforma el cociente de Rayleigh generalizado en uno estándar:

$$\frac{\mathbf{a}^t \mathbf{B} \mathbf{a}}{\mathbf{a}^t \mathbf{W} \mathbf{a}} = \frac{\mathbf{b}^t \mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2} \mathbf{b}}{\mathbf{b}^t \mathbf{b}}.$$

**Teorema 3.1** (Primer factor discriminante canónico). *El máximo de  $\frac{\mathbf{a}^t \mathbf{B} \mathbf{a}}{\mathbf{a}^t \mathbf{W} \mathbf{a}}$  se alcanza con*

$$\mathbf{a}_1 = \mathbf{W}^{-1/2} \mathbf{b}_1,$$

donde  $\mathbf{b}_1$  es el vector propio unitario asociado al mayor valor propio  $\lambda_1$  de la matriz simétrica semidefinida positiva  $\mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2}$ . El primer factor discriminante canónico es

$$\mathbf{y}_1 = \mathbf{H}\mathbf{X}\mathbf{a}_1.$$

*Demostración.* El cociente  $\frac{\mathbf{b}^t \mathbf{M} \mathbf{b}}{\mathbf{b}^t \mathbf{b}}$ , con  $\mathbf{M} = \mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2}$  simétrica, alcanza su máximo  $\lambda_1$  en el vector propio  $\mathbf{b}_1$  asociado a dicho valor propio (cociente de Rayleigh clásico). Deshaciendo el cambio,  $\mathbf{a}_1 = \mathbf{W}^{-1/2} \mathbf{b}_1$ .  $\square$

#### 3.3.1. Conexión con los valores propios de $\mathbf{W}^{-1}\mathbf{B}$

**Proposición 3.2.**  $\mathbf{a}_1 = \mathbf{W}^{-1/2} \mathbf{b}_1$  es vector propio de  $\mathbf{W}^{-1}\mathbf{B}$  asociado al mismo valor propio  $\lambda_1$ , es decir,  $\mathbf{W}^{-1}\mathbf{B}\mathbf{a}_1 = \lambda_1 \mathbf{a}_1$ .

*Demostración.* De  $\mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2} \mathbf{b}_1 = \lambda_1 \mathbf{b}_1$  se tiene, premultiplicando por  $\mathbf{W}^{-1/2}$ :

$$\mathbf{W}^{-1} \mathbf{B} \underbrace{\mathbf{W}^{-1/2} \mathbf{b}_1}_{\mathbf{a}_1} = \lambda_1 \underbrace{\mathbf{W}^{-1/2} \mathbf{b}_1}_{\mathbf{a}_1}. \quad \square$$

En la práctica es más eficiente trabajar con la matriz simétrica  $\mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2}$ , pero conceptualmente el resultado se enuncia equivalentemente en términos de  $\mathbf{W}^{-1}\mathbf{B}$ .

### 3.4. Número de factores discriminantes canónicos

**Teorema 3.3.** Si  $\text{rango}(\mathbf{X}) = p$ , el número máximo de factores discriminantes canónicos es

$$s = \min(p, g - 1).$$

*Demostración.* Como  $\mathbf{W}$  es invertible,  $\text{rango}(\mathbf{W}^{-1}\mathbf{B}) = \text{rango}(\mathbf{B})$ . Ahora bien,

$$\text{rango}(\mathbf{B}) = \text{rango}(\mathbf{X}^t(\mathbf{D} - \mathbf{J})\mathbf{X}) \leq \min(p, \text{rango}(\mathbf{D} - \mathbf{J})) = \min(p, g - 1).$$

La igualdad se alcanza cuando  $\text{rango}(\mathbf{X}) = p$ . La cota  $g - 1$  refleja que conocer las medias de  $g - 1$  grupos determina la del último:  $\bar{\mathbf{x}}_g = (n\bar{\mathbf{x}} - \sum_{k=1}^{g-1} n_k \bar{\mathbf{x}}_k) / n_g$ .  $\square$

### 3.5. Los $s$ factores discriminantes son no correlacionados

**Definición 3.4.** Los factores discriminantes canónicos son las variables  $\mathbf{y}_k = \mathbf{H}\mathbf{X}\mathbf{a}_k$ ,  $k = 1, \dots, s$ , con  $\mathbf{a}_k = \mathbf{W}^{-1/2}\mathbf{b}_k$  y  $\mathbf{b}_k$  el  $k$ -ésimo vector propio de  $\mathbf{W}^{-1/2}\mathbf{B}\mathbf{W}^{-1/2}$  (ordenados por valor propio decreciente).

**Teorema 3.5.**  $\text{Cov}(\mathbf{y}_1, \mathbf{y}_2) = 0$ ; más en general, los  $s$  factores discriminantes canónicos tienen covarianzas dos a dos nulas.

*Demostración.* Como los factores están centrados ( $\mathbf{H}\mathbf{X}\mathbf{a}_k$  tiene media cero),

$$\begin{aligned} n \text{Cov}(\mathbf{y}_1, \mathbf{y}_2) &= \mathbf{a}_1^t \mathbf{X}^t \mathbf{H} \mathbf{X} \mathbf{a}_2 = \mathbf{a}_1^t \mathbf{T} \mathbf{a}_2 = \mathbf{a}_1^t (\mathbf{W} + \mathbf{B}) \mathbf{a}_2 \\ &= \mathbf{a}_1^t \mathbf{W} \mathbf{a}_2 + \mathbf{a}_1^t \mathbf{B} \mathbf{a}_2. \end{aligned}$$

Para el primer sumando, con  $\mathbf{a}_k = \mathbf{W}^{-1/2}\mathbf{b}_k$ :

$$\mathbf{a}_1^t \mathbf{W} \mathbf{a}_2 = \mathbf{b}_1^t \underbrace{\mathbf{W}^{-1/2} \mathbf{W} \mathbf{W}^{-1/2}}_{\mathbf{I}} \mathbf{b}_2 = \mathbf{b}_1^t \mathbf{b}_2 = 0,$$

pues  $\mathbf{b}_1, \mathbf{b}_2$  son vectores propios de una matriz simétrica para valores propios distintos, luego ortogonales. Para el segundo sumando, usando la simetría de  $\mathbf{W}^{-1/2}$ :

$$\mathbf{a}_1^t \mathbf{B} \mathbf{a}_2 = \mathbf{b}_1^t \underbrace{\mathbf{W}^{-1/2} \mathbf{B} \mathbf{W}^{-1/2}}_{\mathbf{M}} \mathbf{b}_2 = \mathbf{b}_1^t (\lambda_2 \mathbf{b}_2) = \lambda_2 \mathbf{b}_1^t \mathbf{b}_2 = 0. \quad \square$$

### 3.6. Interpretación de los factores

Cada factor  $\mathbf{y}_k$  proyecta los datos en la dirección  $\mathbf{a}_k$  que, dentro de la restricción  $\mathbf{a}^t \mathbf{W} \mathbf{a} = 1$ , maximiza  $\mathbf{a}^t \mathbf{B} \mathbf{a}$  sujeto a ser no correlacionado con los factores anteriores. El valor propio  $\lambda_k$  mide la separación entre grupos en esa dirección relativa a la dispersión intragrupos:

$$\lambda_k = \frac{\mathbf{a}_k^t \mathbf{B} \mathbf{a}_k}{\mathbf{a}_k^t \mathbf{W} \mathbf{a}_k}, \quad \text{con } \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_s \geq 0.$$

La proporción de discriminación explicada por el  $k$ -ésimo factor es

$$\rho_k = \frac{\lambda_k}{\sum_{j=1}^s \lambda_j}.$$

## 4. Análisis Discriminante: clasificación

### 4.1. El problema de clasificación

Sea  $\mathbf{x} = (x_1, \dots, x_p)^t$  el vector de características observadas en un individuo cuyo grupo de pertenencia se desconoce. Hay  $g$  subpoblaciones y se dispone de muestras de entrenamiento bien clasificadas. El objetivo es asignar  $\mathbf{x}$  a uno de los  $g$  grupos.

Las situaciones habituales son:

1. La distribución de  $\mathbf{x}$  está *completamente especificada* en cada subpoblación.
2. La distribución es *conocida salvo por el valor de algunos parámetros*.
3. No se conoce la distribución de  $\mathbf{x}$  (métodos no paramétricos).

### 4.2. El método de máxima verosimilitud

**Definición 4.1.** El criterio de máxima verosimilitud asigna el individuo  $\mathbf{x}$  al grupo  $k$  si

$$\mathcal{L}(\mathbf{x} | k) = \max_{j=1, \dots, g} \mathcal{L}(\mathbf{x} | j),$$

donde  $\mathcal{L}(\mathbf{x} | j)$  denota la función de verosimilitud de  $\mathbf{x}$  en la subpoblación  $j$ .

#### 4.2.1. Caso univariante con igual varianza

Para  $g = 2$  con  $X | 1 \sim \mathcal{N}(\mu_1, \sigma)$  y  $X | 2 \sim \mathcal{N}(\mu_2, \sigma)$ , el cociente de verosimilitudes es:

$$\frac{\mathcal{L}(x | 1)}{\mathcal{L}(x | 2)} = \exp\left(-\frac{(x - \mu_1)^2 - (x - \mu_2)^2}{2\sigma^2}\right).$$

**Proposición 4.2.** El individuo  $x$  se clasifica en la subpoblación 1 si y solo si  $x$  está más próximo a  $\mu_1$  que a  $\mu_2$ , es decir,

$$|x - \mu_1| < |x - \mu_2| \iff x > \bar{\mu} := \frac{\mu_1 + \mu_2}{2} \text{ si } \mu_1 > \mu_2 \text{ (y al revés si } \mu_1 < \mu_2).$$

*Demostración.*  $\mathcal{L}(x | 1) > \mathcal{L}(x | 2)$  iff  $-(x - \mu_1)^2 + (x - \mu_2)^2 > 0$  iff  $2x(\mu_1 - \mu_2) > \mu_1^2 - \mu_2^2 = (\mu_1 - \mu_2)(\mu_1 + \mu_2)$ . Dividiendo por  $\mu_1 - \mu_2$  (y cambiando la desigualdad si es negativo) se obtiene el resultado.  $\square$

#### 4.2.2. Caso univariante con desiguales varianzas

Para  $X | j \sim \mathcal{N}(\mu_j, \sigma_j)$ ,  $j = 1, 2$ , la frontera de decisión es *cuadrática*:

$$z_1^2 + 2 \ln \sigma_1 < z_2^2 + 2 \ln \sigma_2, \quad z_j = \frac{x - \mu_j}{\sigma_j}.$$

### 4.3. Caso multivariante: discriminante lineal (LDA)

#### 4.3.1. Dos subpoblaciones

Suponemos  $\mathbf{x} \mid j \sim \mathcal{N}_p(\boldsymbol{\mu}_j, \Sigma)$ ,  $j = 1, 2$ , con la *misma* matriz de covarianzas  $\Sigma$ .

**Teorema 4.3** (Regla lineal de clasificación). *El criterio de máxima verosimilitud clasifica  $\mathbf{x}$  en la subpoblación 1 si*

$$(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \Sigma^{-1} \left( \mathbf{x} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right) > 0.$$

*Demostración.* El cociente de verosimilitudes es

$$\frac{\mathcal{L}(\mathbf{x} \mid 1)}{\mathcal{L}(\mathbf{x} \mid 2)} = \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_1)^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_2)^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \right\}.$$

Expandiendo las formas cuadráticas:

$$\begin{aligned} & - (\mathbf{x} - \boldsymbol{\mu}_1)^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_1) + (\mathbf{x} - \boldsymbol{\mu}_2)^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_2) \\ & = 2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \Sigma^{-1} \mathbf{x} - \boldsymbol{\mu}_1^t \Sigma^{-1} \boldsymbol{\mu}_1 + \boldsymbol{\mu}_2^t \Sigma^{-1} \boldsymbol{\mu}_2 \\ & = 2(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^t \Sigma^{-1} \left( \mathbf{x} - \frac{\boldsymbol{\mu}_1 + \boldsymbol{\mu}_2}{2} \right). \end{aligned}$$

Esta cantidad es positiva iff se cumple la condición del enunciado.  $\square$

#### 4.3.2. $g$ subpoblaciones con igual $\Sigma$

**Teorema 4.4.** *Con  $g$  subpoblaciones gaussianas de igual covarianza  $\Sigma$ , el criterio de máxima verosimilitud equivale a la distancia de Mahalanobis:*

$$\text{Clasificar en } k = \arg \min_{j=1, \dots, g} D^2(\mathbf{x}, \boldsymbol{\mu}_j) := (\mathbf{x} - \boldsymbol{\mu}_j)^t \Sigma^{-1} (\mathbf{x} - \boldsymbol{\mu}_j).$$

*Demostración.* Maximizar  $\mathcal{L}(\mathbf{x} \mid j) = |2\pi\Sigma|^{-1/2} \exp\{-\frac{1}{2}D^2(\mathbf{x}, \boldsymbol{\mu}_j)\}$  sobre  $j$  equivale a minimizar  $D^2(\mathbf{x}, \boldsymbol{\mu}_j)$  sobre  $j$ , dado que el factor  $|2\pi\Sigma|^{-1/2}$  es constante en  $j$ .  $\square$

Expandiendo  $D^2(\mathbf{x}, \boldsymbol{\mu}_j)$  y omitiendo el término  $\mathbf{x}^t \Sigma^{-1} \mathbf{x}$  (constante en  $j$ ), la regla equivale a maximizar las **funciones discriminantes lineales**:

$$\delta_j(\mathbf{x}) = \boldsymbol{\mu}_j^t \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_j^t \Sigma^{-1} \boldsymbol{\mu}_j, \quad j = 1, \dots, g.$$

### 4.4. Discriminante cuadrático (QDA)

Cuando las matrices de covarianzas  $\Sigma_j$  son distintas entre grupos, el criterio de máxima verosimilitud ya no produce fronteras lineales.

**Teorema 4.5.** *Con  $\mathbf{x} \mid j \sim \mathcal{N}_p(\boldsymbol{\mu}_j, \Sigma_j)$  y matrices  $\Sigma_j$  posiblemente distintas, se clasifica  $\mathbf{x}$  en el grupo  $k$  donde se maximiza la **función discriminante cuadrática**:*

$$q_j(\mathbf{x}) = -\frac{1}{2} \mathbf{x}^t \Sigma_j^{-1} \mathbf{x} + \boldsymbol{\mu}_j^t \Sigma_j^{-1} \mathbf{x} - \frac{1}{2} \ln |\Sigma_j| - \frac{1}{2} \boldsymbol{\mu}_j^t \Sigma_j^{-1} \boldsymbol{\mu}_j.$$

*Demostración.* Tomando logaritmos en la verosimilitud y descartando la constante  $-\frac{p}{2} \ln(2\pi)$ :

$$\ln \mathcal{L}(\mathbf{x} \mid j) = -\frac{1}{2} \ln |\Sigma_j| - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}_j)^t \Sigma_j^{-1} (\mathbf{x} - \boldsymbol{\mu}_j).$$

Expandiendo la forma cuadrática se obtiene directamente  $q_j(\mathbf{x})$ .  $\square$

Las fronteras de decisión son ahora superficies cuadráticas (hiperboloides, elipsoides o paraboloides según los  $\Sigma_j$ ), en contraste con los hiperplanos del caso lineal.

## 4.5. El método bayesiano

En muchas aplicaciones existe información a priori sobre la frecuencia relativa de cada grupo en la población.

**Definición 4.6.** Sean  $\pi_j > 0$ ,  $\sum_j \pi_j = 1$ , las **probabilidades a priori** de pertenencia al grupo  $j$ . El **método bayesiano** asigna  $\mathbf{x}$  al grupo  $k$  donde se maximiza la probabilidad a posteriori:

$$\max_{j=1,\dots,g} \mathcal{L}(\mathbf{x} | j) \pi_j.$$

**Observación 4.7.** La maximización de  $\mathcal{L}(\mathbf{x} | j)\pi_j$  es equivalente a la de  $\mathcal{L}(j | \mathbf{x}) = \mathcal{L}(\mathbf{x} | j)\pi_j / \mathcal{L}(\mathbf{x})$ , pues el denominador no depende de  $j$ . Si  $\pi_j = 1/g$  para todo  $j$ , el método bayesiano se reduce al de máxima verosimilitud.

### 4.5.1. LDA bayesiano

Para subpoblaciones gaussianas con igual  $\Sigma$ , el método bayesiano maximiza sobre  $j$ :

$$\delta_j^B(\mathbf{x}) = \boldsymbol{\mu}_j^t \Sigma^{-1} \mathbf{x} - \frac{1}{2} \boldsymbol{\mu}_j^t \Sigma^{-1} \boldsymbol{\mu}_j + \ln \pi_j.$$

### 4.5.2. QDA bayesiano

Con matrices  $\Sigma_j$  distintas, se maximiza:

$$q_j^B(\mathbf{x}) = q_j(\mathbf{x}) + \ln \pi_j.$$

## 5. Estimación de parámetros desconocidos

En la práctica, los parámetros  $\boldsymbol{\mu}_j$  y  $\Sigma$  (o  $\Sigma_j$ ) son desconocidos y se sustituyen por sus estimadores de máxima verosimilitud a partir de las muestras de entrenamiento:

$$\hat{\boldsymbol{\mu}}_j = \bar{\mathbf{x}}_j, \quad \hat{\Sigma} = \frac{1}{n} \mathbf{W} = \frac{1}{n} \sum_{k=1}^g n_k \mathbf{S}_k \quad (\text{caso LDA}),$$

$$\hat{\Sigma}_j = \mathbf{S}_j \quad (\text{caso QDA}).$$

Las reglas de clasificación resultantes se denominan, respectivamente, **LDA muestral** y **QDA muestral**.

**Observación 5.1.** Cuando  $p$  es grande respecto a  $n$ , la estimación de  $\Sigma$  es inestable. En esos casos es habitual usar versiones regularizadas:  $\hat{\Sigma}_\alpha = (1 - \alpha)\hat{\Sigma} + \alpha\mathbf{I}$ , con  $\alpha \in [0, 1]$  escogido por validación cruzada.

## 6. Resumen y relación entre métodos

Método	Supuesto	Frontera	Criterio
LDA	$\Sigma_j = \Sigma$	Lineal	$\max_j \delta_j(\mathbf{x})$
LDA bayesiano	$\Sigma_j = \Sigma$ , $\pi_j$ conocidas	Lineal	$\max_j \delta_j^B(\mathbf{x})$
QDA	$\Sigma_j$ distintas	Cuadrática	$\max_j q_j(\mathbf{x})$
QDA bayesiano	$\Sigma_j$ distintas, $\pi_j$ conocidas	Cuadrática	$\max_j q_j^B(\mathbf{x})$
Análisis factorial	Descriptivo	—	Vectores propios de $\mathbf{W}^{-1}\mathbf{B}$

La conexión entre los dos grandes enfoques es la siguiente: el **análisis factorial discriminante** (aspecto descriptivo) produce las direcciones  $\mathbf{a}_1, \dots, \mathbf{a}_s$  que maximizan progresivamente la separación entre grupos; el **análisis discriminante clásico** (aspecto decisional) usa reglas probabilísticas para asignar nuevos individuos. Ambos comparten la descomposición  $\mathbf{T} = \mathbf{W} + \mathbf{B}$  como herramienta algebraica fundamental.