

# 1 Análisis Multivariado - Práctica 5

## 1.1 Clasificación y Coordenadas discriminantes

1. Sea  $\mathbf{x} \sim Bi(n, \theta_i)$  en  $G_i$ , con  $i = 1, 2$ . Encontrar la regla óptima de clasificación y mostrar que puede llevarse a una función discriminante lineal.
2. Supongamos que  $\mathbf{x} \sim Exp(\lambda_i)$  en  $G_i$ , con  $i = 1, 2$ .
  - (a) Encontrar la regla óptima de clasificación y expresarla como una función discriminante lineal.
  - (b) Calcular la probabilidad total de mala clasificación  $P(\mathcal{R}, \mathbf{f})$  cuando  $\pi_1 = \pi_2$ .
  - (c) Tomemos  $\lambda_1 = 1$  y  $\lambda_2 = \lambda > 1$ . Estudiar el límite de  $P(\mathcal{R}, \mathbf{f})$  cuando  $\lambda \rightarrow \infty$ . Sacar conclusiones.

3. Una regla de clasificación es *minimax* si las regiones que definen el criterio de clasificación se buscan de modo que minimicen  $\max\{P(1 | 2), P(2 | 1)\}$ .

- (a) Dado  $\alpha \in (0, 1)$ , verificar que

$$\max\{P(1 | 2), P(2 | 1)\} \geq (1 - \alpha)P(1 | 2) + \alpha P(2 | 1)$$

- (b) Para cada  $\alpha$ , encontrar la regla que minimiza el lado derecho de la ecuación anterior.
- (c) Probar que la regla minimax está dada por

$$R_1 = \left\{ x : \frac{f_1(x)}{f_2(x)} > c \right\}$$

donde  $c$  satisface que  $P(1 | 2) = P(2 | 1)$ .

4. Sea  $\mathbf{x}$  un vector aleatorio, y sean  $\mu_1 = E(\mathbf{x} | G_1)$ ,  $\mu_2 = E(\mathbf{x} | G_2)$ . Supongamos que la matriz de covarianza  $\Sigma = E((\mathbf{x} - \mu_i)(\mathbf{x} - \mu_i)')$  ( $i = 1, 2$ ) es la misma para ambas poblaciones, que los costos son iguales y que las  $\pi_i$  también lo son.
  - (a) Si  $B$  se define como  $B = c(\mu_1 - \mu_2)(\mu_1 - \mu_2)'$  para alguna constante  $c$ , verificar que  $\mathbf{e} = c\Sigma^{-1}(\mu_1 - \mu_2)$  es de hecho un autovector (sin escalar) de  $\Sigma^{-1}B$ .
  - (b) Deducir que cuando  $k = 2$ , el criterio de clasificación de Fisher (es decir el análisis discriminante) coincide con la regla óptima para la normal.
5. Sea el siguiente modelo de regresión:  $z = b'\mathbf{x} + \mathbf{e}$ ,  $\mathbf{e}$  con distribución logística. Supongamos que sólo observamos si  $z > 0$  ó  $z < 0$ , (no su valor). Mostrar que estimar a  $b$  coincide con encontrar el estimador de los parámetros del modelo de clasificación basada en el modelo de regresión logística en el caso  $k = 2$ .

6. Consideremos los datos “iris” del R. Es un conjunto de datos analizados por Fisher que consisten en 4 mediciones realizadas en 50 flores iris de cada una de 3 especies distintas (Setosa, Versicolor y Virginica). Las 4 variables, medidas en centímetros, son

$X_1 =$  Longitud de los sépalos (sepal length)

$X_2 =$  Ancho de los sépalos (sepal width)

$X_3 =$  Longitud de los pétalos (petal length)

$X_4 =$  Ancho de los pétalos (petal width)

- (a) Suponiendo normalidad en los datos, hacer una clasificación lineal y calcular el error de clasificación por los 3 métodos:
- “Ingenuo” (calcular la proporción de datos mal clasificados) o error aparente
  - Validación cruzada
- (b) Ídem b) pero con la clasificación cuadrática y comparar los resultados.

7. Del conjunto de datos “iris” consideremos las variables  $X_2 =$  Ancho de los sépalos y  $X_4 =$  Ancho de los pétalos para las 3 especies de flores.

- (a) Suponiendo que la distribución es normal bivariada para cada población, construir la regla de clasificación cuadrática, asumiendo costos de mala clasificación iguales y probabilidad a priori de pertenecer a cada grupo iguales. Usando esta regla de clasificación recién construida clasificar la nueva observación  $\mathbf{x}_0 = (3.5, 1.75)'$  como perteneciente a alguno de los 3 grupos.
- (b) Supongamos que las matrices de covarianza  $\Sigma_i$  son las mismas para las 3 poblaciones normales bivariadas. Construir la regla de clasificación lineal, asumiendo costos de mala clasificación iguales y probabilidad a priori de pertenecer a cada grupo iguales, y usarla para clasificar la nueva observación  $\mathbf{x}_0 = (3.5, 1.75)'$  como perteneciente a alguno de los 3 grupos. Comparar los resultados obtenidos en b) y c). ¿Cuál enfoque es preferible en este caso?
- (c) Graficar en un scatterplot las regiones halladas en b).
- (d) Usando la clasificación lineal realizada en b), clasificar las observaciones de la muestra. Calcular el error aparente total y la estimación insesgada del error que se obtiene por validación cruzada.

8. Aproximadamente 2 años antes de la bancarrota de algunas empresas se recolectan datos financieros de las mismas, y también se recolectan datos de empresas sanas financieramente alrededor del mismo momento. A continuación figuran los datos de 4 variables:

$X_1 =$  (flujo de caja)/(deuda total)

$X_2 =$  (ingreso neto)/(total de activos)

$X_3 =$  (activos corrientes)/(pasivos corrientes)

$X_4 =$  (activos corrientes)/(ventas netas)

- (a) Graficar los datos para los pares de observaciones  $(X_1, X_2)$ ,  $(X_1, X_3)$  y  $(X_1, X_4)$ . Para alguno de estos pares de variables, ¿tienen aspecto de provenir de una distribución normal bivariada?

- (b) Usando los  $n_1 = 21$  pares de observaciones  $(X_1, X_2)$  de empresas en bancarrota y los  $n_2 = 25$  pares de observaciones  $(X_1, X_2)$  de empresas sanas financieramente, calcular los vectores de medias muestrales  $\bar{X}_1$  y  $\bar{X}_2$  y las matrices de covarianza muestrales  $S_1$  y  $S_2$ .
  - (c) Usando los resultados de b) y asumiendo que las dos muestras aleatorias provienen de dos poblaciones normales, construir la regla de clasificación cuadrática asumiendo  $\pi_1 = \pi_2$  y  $c(1 | 2) = c(2 | 1)$ .
  - (d) Evaluar la performance de la regla de clasificación desarrollada en c) calculando el error aparente total y la estimación del error actual esperado que se obtiene por validación cruzada.
  - (e) Repetir los items c) y d) tomando  $\pi_1 = 0.05$  y  $\pi_2 = 0.95$  y  $c(1 | 2) = c(2 | 1)$ . ¿Es razonable esta elección de probabilidades a priori?
  - (f) Repetir los items b) a e) usando las variables  $(X_1, X_3)$  y  $(X_1, X_4)$ . ¿Parecen ser algunas variables mejores clasificadoras que otras?
- h) Repetir los items b) a e) usando las 4 variables.