

Análisis Multivariado 2 - Práctica 2

Correlación canónica

1. (a) Consideremos un vector $\mathbf{z} \in \mathbb{R}^d$ tal que $\mathcal{E}[\mathbf{z}] = \mathbf{0}$ y $\mathcal{D}[\mathbf{z}] = \Sigma$. Para medir la asociación lineal entre la primera componente z_1 y las demás, $\mathbf{y} = (z_2, \dots, z_d)'$, se define el coeficiente de correlación múltiple al cuadrado $\rho_{1(23\dots d)}^2$ como la mayor correlación (al cuadrado) entre z_1 y cualquier combinación lineal de \mathbf{y} . Es decir,

$$\rho_{1(23\dots d)}^2 = \max_{\beta} \frac{[\text{cov}(z_1, \beta' \mathbf{y})]^2}{\text{var}(z_1) \text{var}(\beta' \mathbf{y})}. \quad (1)$$

Probar que

$$\rho_{1(23\dots d)}^2 = \frac{\sigma_{12} \Sigma_{22}^{-1} \sigma_{21}}{\sigma_{11}}$$

con los parámetros que vienen de la partición

$$\Sigma = \begin{pmatrix} \sigma_{11} & \sigma_{12} \\ \sigma_{21} & \Sigma_{22} \end{pmatrix}.$$

Además, probar que el máximo (1) se realiza en $\beta = \Sigma_{22}^{-1} \sigma_{21}$.

- (b) Supongamos ahora que queremos predecir z_1 mediante una combinación lineal de \mathbf{y} . Entonces se busca

$$\beta^* = \arg \min_{\beta} E \left[(z_1 - \beta' \mathbf{y})^2 \right].$$

Probar que nuevamente se obtiene $\beta^* = \Sigma_{22}^{-1} \sigma_{21}$.

2. Sea $\mathbf{z} = (\mathbf{x}', \mathbf{y}')'$ con $\mathbf{x} \in \mathbb{R}^{d_1}$ e $\mathbf{y} \in \mathbb{R}^{d_2}$. Si $\Sigma = \mathcal{D}[\mathbf{z}]$ es definida positiva, probar que la primera correlación canónica ρ_1 es estrictamente menor que 1.

SUGERENCIA: Usar A3.2 de Seber.

3. Probar que las correlaciones canónicas son invariantes por transformaciones afines. Es decir, las correlaciones canónicas entre \mathbf{x} e \mathbf{y} son las mismas que entre $A\mathbf{x}$ y $B\mathbf{y}$ si A y B son matrices inversibles.

4. Dadas dos variables canónicas u_i y v_j con $i \neq j$, demostrar que $\text{cov}[u_i, v_j] = 0$.

SUGERENCIA: Primero mostrar que $\mathbf{a}_i = \Sigma_{11}^{-1} \Sigma_{12} \mathbf{b}_i$.

5. Usando multiplicadores de Lagrange, probar que ρ_1^2 es el máximo de $(\alpha' \Sigma_{12} \beta)^2$ sujeto a $\alpha' \Sigma_{11} \alpha = 1$ y $\beta' \Sigma_{22} \beta = 1$.

SUGERENCIA: Usar A8.1 de Seber.

6. Sea $\mathbf{z} = (\mathbf{x}', \mathbf{y}')'$ con $\mathbf{x} \in \mathbb{R}^2$ e $\mathbf{y} \in \mathbb{R}^2$ y supongamos que

$$\mathcal{D}[\mathbf{z}] = \sigma^2 \begin{pmatrix} 1 & a & b & b \\ a & 1 & b & b \\ b & b & 1 & c \\ b & b & c & 1 \end{pmatrix}$$

donde a , b y c tienen módulo menor que 1. Encontrar la primera correlación canónica y las correspondientes variables canónicas.

- (a) Hallar las correlaciones canónicas muestrales.
- (b) Determinar el primer par canónico \hat{u}_1 y \hat{v}_1 e interpretar estas cantidades.

9. En un estudio realizado por Waugh (1942) basado en $n = 138$ muestras de trigo de variedad Canadian Hard Red Spring y de la harina hecha a partir de él, se midieron las siguientes variables (5 son mediciones en forma estandarizada para el trigo X_i y 4 para la harina Y_i):

X_1 = estructura del núcleo

X_2 = peso

X_3 = núcleos dañados

X_4 = material extraño

X_5 = proteínas

Y_1 = trigo por barril de harina

Y_2 = ceniza presente en la harina

Y_3 = proteínas en la harina

Y_4 = índice de calidad del gluten

La matriz de correlación fue:

$$R = \left[\begin{array}{cc|cc} R_{11} & R_{12} \\ R_{21} & R_{22} \end{array} \right] = \left[\begin{array}{ccccc|cccc} 1 & & & & & & & & & & \\ 0.754 & 1 & & & & & & & & & \\ -0.69 & -0.712 & 1 & & & & & & & & \\ -0.446 & -0.515 & 0.323 & 1 & & & & & & & \\ 0.692 & 0.412 & -0.444 & -0.334 & 1 & & & & & & \\ \hline -0.605 & -0.722 & 0.737 & 0.527 & -0.383 & 1 & & & & & \\ -0.479 & -0.419 & 0.361 & 0.461 & -0.505 & 0.251 & 1 & & & & \\ -0.780 & 0.542 & -0.546 & -0.393 & 0.737 & -0.490 & -0.434 & 1 & & & \\ -0.152 & -0.102 & 0.172 & -0.019 & -0.148 & 0.250 & -0.079 & -0.163 & 1 & & \end{array} \right]$$

- (a) Hallar las variables canónicas muestrales correspondientes a las correlaciones canónicas.
- (b) Interpretar las variables canónicas muestrales \hat{u}_1 y \hat{v}_1 . ¿Representan en algún sentido la calidad del trigo y la harina, respectivamente?
- (c) ¿Qué proporción de la varianza muestral total de las variables X queda explicada por \hat{u}_1 ? ¿Qué proporción de la varianza muestral total de las variables Y queda explicada por \hat{v}_1 ?

10. En el archivo hogares se encuentran los datos correspondientes a 75 observaciones de hogares españoles. Las primeras cinco variables son el gasto del hogar en distintas partidas y las cuatro siguientes la estructura del hogar.

Y_1 = gastos en alimentación, bebidas y tabaco.

Y_2 = gastos en vestido y calzado.

Y_3 = gastos en menaje.

Y_4 = gastos en transporte y comunicaciones.

Y_5 = gastos en espacimientoy enseñanza.

X_1 = número de personas en la unidad de gastos.

X_2 = número de personas mayores de 14 años.

X_3 = el nivel académico del sustentador principal medido en una escala de 1 a 8.
 X_4 = número de perceptores con ingresos.

- (a) Hallar las correlaciones canónicas muestrales.
- (b) Determinar el primer par canónico \hat{u}_1 y \hat{v}_1 e interpretar estas cantidades.