

# 1 Análisis Multivariado - Práctica 4

## 1.1 Correlación canónica

- (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  $\rho_1$  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  $\rho_1^2$  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.

7. Sea  $\mathbf{z} = (\mathbf{x}', \mathbf{y}')$  con  $\mathbf{x} \in \mathbb{R}^{d_1}$  e  $\mathbf{y} \in \mathbb{R}^{d_2}$ . Supongamos que  $\mathcal{E}[\mathbf{z}] = \mathbf{0}$  y consideramos  $\Sigma = \mathcal{D}[\mathbf{z}]$  con la partición

$$\Sigma = \begin{pmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{pmatrix} = \begin{pmatrix} \mathcal{E}[\mathbf{xx}'] & \mathcal{E}[\mathbf{xy}'] \\ \mathcal{E}[\mathbf{yx}'] & \mathcal{E}[\mathbf{yy}'] \end{pmatrix}.$$

Se quiere predecir  $\mathbf{y}$  mediante  $k$  combinaciones lineales no correlacionadas de  $\mathbf{x}$ . Es decir, si definimos  $\mathbf{u} = A\mathbf{x}$  con  $A$  de  $k \times d_1$  tal que  $\mathcal{D}[\mathbf{u}] = \Sigma_{\mathbf{u}} = A\Sigma_{11}A' = I_k$ , queremos encontrar un predictor lineal de  $\mathbf{y}$  de la forma  $B\mathbf{u}$  (donde  $B$  será de  $d_2 \times k$ ). El criterio será elegir  $A$  y  $B$  que minimicen  $E[\|\mathbf{y} - B\mathbf{u}\|^2]$ .

- (a) Probar que, fijada la matriz  $A$ , se tiene que

$$E[\|\mathbf{y} - B\mathbf{u}\|^2] \geq E[\|\mathbf{y} - B^*\mathbf{u}\|^2]$$

para  $B^* = \Sigma_{21}A'(A\Sigma_{11}A')^{-1}$ . Mostrar además que

$$E[\|\mathbf{y} - B^*\mathbf{u}\|^2] = \text{tr}(\Sigma_{22}) - \text{tr}\left(\Sigma_{21}A'(A\Sigma_{11}A')^{-1}A\Sigma_{12}\right),$$

por lo que el problema se reduce a encontrar

$$A^* = \arg \max_{A\Sigma_{11}A'=I_k} \text{tr}\left(\Sigma_{21}A'(A\Sigma_{11}A')^{-1}A\Sigma_{12}\right). \quad (2)$$

- (b) Mostrar que la matriz  $A^*$  que cumple (2) tiene como filas a los autovectores correspondientes a los primeros  $k$  autovalores de  $\Sigma_{11}^{-1}\Sigma_{12}\Sigma_{21}$  y tales que son ortogonales con la distancia dada por  $\Sigma_{11}$ .

SUGERENCIA: Usar el lema que sigue: Sea  $Q \in \mathbb{R}^{d \times d}$ ,  $Q \geq 0$  y  $k < d$ ,  $C \in \mathbb{R}^{k \times d}$  tal que  $CC' = I_k$  entonces

$$\sum_{i=1}^k \lambda_i(CQC') \leq \sum_{i=1}^k \lambda_i(Q)$$

donde  $\lambda_1(Q) \geq \dots \geq \lambda_d(Q)$  son los autovalores de  $Q$  y  $\lambda_1(CQC') \geq \dots \geq \lambda_k(CQC')$  los autovalores de  $CQC'$ .

8. En un estudio de pobreza, crimen y disuasión, Parker y Smith reportaron ciertos resúmenes de estadísticas criminales en varios estados de EEUU para los años 1970 y 1973. Una parte de la matriz de correlación muestral aparece más abajo. Las variables son:

$X_1$  = homicidios no primarios cometidos durante 1970

$X_2$  = homicidios primarios (homicidios que involucran parientes o conocidos) cometidos durante 1970

$Y_1$  = severidad de los castigos (mediana de los meses cumplidos en prisión) en 1970

$Y_2$  = certeza del castigo (número de admisiones a prisión dividido por el número de homicidios) en 1970

$$R = \left[ \begin{array}{cc|cc} & & & & & \\ & & & & & \\ R_{11} & R_{12} & & & & \\ R_{21} & R_{22} & & & & \\ & & & & & \end{array} \right] = \left[ \begin{array}{cc|cc} 1 & 0.615 & -0.111 & -0.266 \\ 0.615 & 1 & -0.195 & -0.085 \\ -0.111 & -0.195 & 1 & -0.269 \\ -0.266 & -0.085 & -0.269 & 1 \end{array} \right]$$

