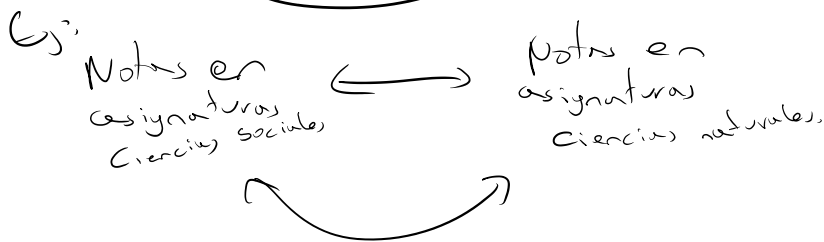
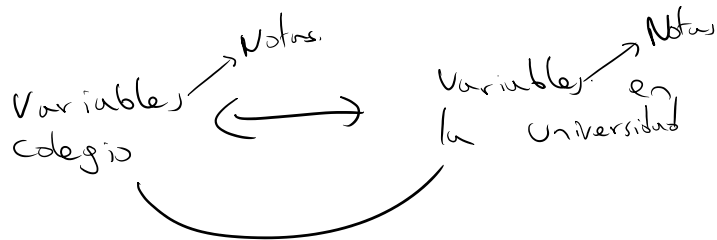


Buscar la asociación de conjuntos de variables.

Ej: Evaluar la relación entre el desempeño escolar y desempeño universitario de alumnos de MACC.



## Metodología general

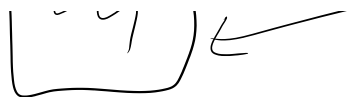
Estudiar la correlación entre comb. lineales de variables del grupo 1 con comb. lineales de variables del grupo 2.



- Buscamos la combinación lineal con correlación más alta.
- Luego, buscamos la segunda comb. lineal no correlacionada con la anterior con correlación más alta.

Idealmente buscamos reducir información de muchas dimensiones de relaciones entre variables de dos grupos a pocos pares de variantes canónicas.





$$\rho = \begin{pmatrix} \rho_{11} & \rho_{12} \\ \rho_{21} & \rho_{22} \end{pmatrix}$$

Supongamos que el primer grupo tiene  $p$  variables  
y lo representamos con el vector  $X^{(1)}$

Supongamos que el segundo grupo tiene  $q$  variables  
y lo representamos con  $X^{(2)}$

Supondremos siempre  $p \leq q$

Definimos

$$E(X^{(1)}) = \mu^{(1)}$$

$$\text{cov}(X^{(1)}) = \Sigma_{11}$$

$$E(X^{(2)}) = \mu^{(2)}$$

$$\text{cov}(X^{(2)}) = \Sigma_{22}$$

$$\text{cov}(X^{(1)}, X^{(2)}) = \Sigma_{12} = \Sigma_{21}'$$

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

$$\rho = \begin{pmatrix} \rho_{11} & \rho_{12} \\ \rho_{21}' & \rho_{22} \end{pmatrix}$$

Obs: Las  $p \times q$  asociaciones lineales entre  $X^{(1)}$  y  $X^{(2)}$   
están en  $\Sigma_{12}$

Sean

$$U = a' X^{(1)}, \quad V = b' X^{(2)}$$

$a$  y  $b$  son vectores  
de coeficientes

$$\text{Var}(U) = a' \Sigma_{11} a$$

$$\text{Var}(U) = a' \Sigma_{11} a$$

$$\text{Var}(V) = b' \Sigma_{22} b$$

$$\text{Cov}(U, V) = a' \Sigma_{12} b$$

→ TAREA.

Buscar  $a, b$  tal que

$$\text{cor}(U, V) = \frac{a' \Sigma_{12} b}{\sqrt{a' \Sigma_{11} a} \sqrt{b' \Sigma_{22} b}} \quad (\text{pepe})$$

Sea máxima.

i) El primer par de variantes canónicas son las comb. lineales  $U_1, V_1$  con varianzas 1 que maximiza  $\alpha$  (pepe)

ii) El segundo par de variantes canónicas son comb. lineales  $U_2, V_2$  con varianzas 1, que maximiza  $\alpha$  (pepe) y no están correlacionadas con  $U_1$  y  $V_1$

iii) El  $k$ -ésimo par de variantes canónicas son las comb. lineales  $U_k, V_k$  que maximizan  $\alpha$  (pepe) y no están correlacionadas con los  $k-1$  pares anteriores

Suponga  $p \leq q$  y  $X^{(1)}, X^{(2)}$  como lo definimos anteriormente +  $q$ .  $\Sigma = \text{cov}(X)$

$$U = a' X^{(1)} \quad V = b' X^{(2)}$$

o. 75

$$S \quad U = a' X \quad V = b' X$$

Entonces

$$1) \max_{a, b} \text{cor}(U, V) = \rho^* \rightarrow 0.75$$

se obtiene con  $U_1 = e_1' \sum_{11}^{-1/2} X^{(1)}$   
 $V_1 = f_1' \sum_{22}^{-1/2} X^{(2)}$

El  $K$ -ésimo par de variables canónicas  $K=2, 3, \dots, p$

$$U_K = e_K' \sum_{11}^{-1/2} X^{(1)} \quad V_K = f_K' \sum_{22}^{-1/2} X^{(2)}$$

no correlacionados con los  $K-1$  pares anteriores,

donde  $\rho_1^* \geq \rho_2^* \geq \dots \geq \rho_p^*$  son los valores propios de la matriz

$$\sum_{11}^{-1/2} \sum_{12} \sum_{22}^{-1} \sum_{21} \sum_{11}^{-1/2}$$

Con vectores propios asociados  $e_1, e_2, \dots, e_p$ .

Los  $f_1, \dots, f_p$  son los vectores propios de la matriz

$$\sum_{22}^{-1/2} \sum_{21} \sum_{11}^{-1} \sum_{12} \sum_{22}^{-1/2}$$

Estas variables satisfacen

$$\text{var}(U_K) = \text{var}(V_K) = 1 \quad \forall K$$

$$\text{cov}(U_K, U_l) = \text{cov}(U_K, V_l) = 0 \quad \forall K \neq l$$

$$\text{cov}(V_K, V_l) = \text{cov}(V_K, U_l) = 0 \quad \forall K \neq l$$

$$\text{cov}(U_k, V_l) = \text{cov}(U_k, V_l) = 0$$

$$\forall k \neq l$$

**Example 10.1 (Calculating canonical variates and canonical correlations for standardized variables)** Suppose  $\mathbf{Z}^{(1)} = [Z_1^{(1)}, Z_2^{(1)}]'$  are standardized variables and  $\mathbf{Z}^{(2)} = [Z_1^{(2)}, Z_2^{(2)}]'$  are also standardized variables. Let  $\mathbf{Z} = [\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}]'$  and

$$\text{Cov}(\mathbf{Z}) = \begin{bmatrix} \rho_{11} & \rho_{12} \\ \rho_{21} & \rho_{22} \end{bmatrix} = \begin{bmatrix} 1.0 & .4 & .5 & .6 \\ .4 & 1.0 & .3 & .4 \\ .5 & .3 & 1.0 & .2 \\ .6 & .4 & .2 & 1.0 \end{bmatrix}$$

Then

$$\rho_{11}^{-1/2} = \begin{bmatrix} 1.0681 & -.2229 \\ -.2229 & 1.0681 \end{bmatrix}$$

$$\rho_{22}^{-1} = \begin{bmatrix} 1.0417 & -.2083 \\ -.2083 & 1.0417 \end{bmatrix}$$

and

$$\rho_{11}^{-1/2} \rho_{12} \rho_{22}^{-1} \rho_{21} \rho_{11}^{-1/2} = \begin{bmatrix} .4371 & .2178 \\ .2178 & .1096 \end{bmatrix}$$

The eigenvalues  $\rho_1^*$  and  $\rho_2^*$  of  $\rho_{11}^{-1/2} \rho_{12} \rho_{22}^{-1} \rho_{21} \rho_{11}^{-1/2}$  are obtained from

$$0 = \begin{vmatrix} .4371 - \lambda & .2178 \\ .2178 & .1096 - \lambda \end{vmatrix} = (.4371 - \lambda)(.1096 - \lambda) - (.2178)^2$$

$$= \lambda^2 - .5467\lambda + .0005$$

er 10 Canonical Correlation Analysis

yielding  $\rho_1^* = .5458$  and  $\rho_2^* = .0009$ . The eigenvector  $\mathbf{e}_1$  follows from the vector equation

$$\begin{bmatrix} .4371 & .2178 \\ .2178 & .1096 \end{bmatrix} \mathbf{e}_1 = (.5458) \mathbf{e}_1$$

Thus,  $\mathbf{e}_1' = [.8947, .4466]$  and

$$\mathbf{a}_1 = \rho_{11}^{-1/2} \mathbf{e}_1 = \begin{bmatrix} .8561 \\ .2776 \end{bmatrix}$$

From Result 10.1,  $\mathbf{f}_1 \propto \rho_{22}^{-1/2} \rho_{21} \rho_{11}^{-1/2} \mathbf{e}_1$  and  $\mathbf{b}_1 = \rho_{22}^{-1/2} \mathbf{f}_1$ . Consequently,

$$\mathbf{b}_1 \propto \rho_{22}^{-1/2} \rho_{21} \mathbf{a}_1 = \begin{bmatrix} .3959 & .2292 \\ .5209 & .3547 \end{bmatrix} \begin{bmatrix} .8561 \\ .2776 \end{bmatrix} = \begin{bmatrix} .4026 \\ .5443 \end{bmatrix}$$

We must scale  $\mathbf{b}_1$  so that

$$\text{Var}(V_1) = \text{Var}(\mathbf{b}_1' \mathbf{Z}^{(2)}) = \mathbf{b}_1' \rho_{22} \mathbf{b}_1 = 1$$

The vector  $[\mathbf{b}_1]'$  gives

$$\begin{bmatrix} .4026 & .5443 \end{bmatrix} \begin{bmatrix} .2 & 1.0 \\ 1.0 & .2 \end{bmatrix} \begin{bmatrix} .4026 \\ .5443 \end{bmatrix} = .5460$$

Using  $\sqrt{.5460} = .7389$ , we take

$$\mathbf{b}_1 = \frac{1}{.7389} \begin{bmatrix} .4026 \\ .5443 \end{bmatrix} = \begin{bmatrix} .5448 \\ .7366 \end{bmatrix}$$

The first pair of canonical variates is

$$U_1 = \mathbf{a}_1' \mathbf{Z}^{(1)} = .86Z_1^{(1)} + .28Z_2^{(1)}$$

$$V_1 = \mathbf{b}_1' \mathbf{Z}^{(2)} = .54Z_1^{(2)} + .74Z_2^{(2)}$$

and their canonical correlation is

$$\rho_1^* = \sqrt{\rho_1^*} = \sqrt{.5458} = .74$$

$$\rho_1^* = 0.73$$

$$\rho_2^* = 0.3$$