

# Relaciones entre dos conjuntos de variables

Básicamente se tienen dos situaciones cuando queremos relacionar dos conjuntos de variables:

- 1 **Situación simétrica:** se desea tratar ambos grupos de variables del mismo modo, donde no se le da ninguna preferencia a ninguno de los dos conjuntos de variables para explicar al otro. Se busca investigar globalmente la relación entre ambos conjuntos de variables.
- 2 **Situación asimétrica:** corresponde al caso donde un conjunto de variables explican a las variables del segundo conjunto y no al revés.

# Relaciones entre dos conjuntos de variables

- **Ejemplo de la situación simétrica:** cuando se considera que el primer grupo de variables es de rendimiento escolar y el segundo de uso del tiempo de ocio entre los estudiante. No existe claramente un conjunto de variables que sea la causa del otro conjunto.
- **Ejemplo del la situación asimétrica:** el primer grupo de variables mide el rendimiento en secundaria y el segundo en la universidad. Claramente las primeras pueden causar las segundas pero no al revés. Los modelos de regresion multivariada son una herramienta para analizar los datos asimétricos.

La situación simétrica la estudio inicialmente Hotelling en 1936, como una extensión de los componentes principales y se denomina **análisis de correlación canónica**

# Análisis de Correlación Canónica

- El análisis de correlación canónica busca identificar y cuantificar las asociaciones entre dos conjuntos de variables.
- El análisis de correlación canónica se enfoca en la correlación entre una combinación lineal de las variables en un conjunto y una combinación lineal de las variables en otro conjunto.
- La idea es determinar primero el par de combinaciones lineales que tienen la mayor correlación. A continuación, se determina un segundo par de combinaciones lineales que tienen la mayor correlación entre todos los pares no correlacionados con el primer par seleccionado, y así sucesivamente.
- Los pares de combinaciones lineales se denominan **variables canónicas**, y sus correlaciones se llaman **correlaciones canónicas**.

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

- Las **correlaciones canónicas** miden la fuerza de asociación entre los dos conjuntos de variables.
- El aspecto de maximización de la técnica representa un intento de concentrar una relación de alta dimensionalidad entre dos conjuntos de variables en unos cuantos pares de **variables canónicas**
- El primer grupo, de  $p$  variables, se representa por el vector aleatorio  $\mathbf{x}^{(1)}$ ,  $(p \times 1)$ . El segundo grupo, de  $q$  variables, se representa por el vector aleatorio  $\mathbf{x}^{(2)}$ ,  $(q \times 1)$ .
- Se asumirá en el desarrollo teórico, que  $\mathbf{x}^{(1)}$  representa el conjunto *más pequeño*, de modo que  $p \leq q$ .

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

Para los vectores aleatorios  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$ , sea

$$E(\mathbf{x}^{(1)}) = \boldsymbol{\mu}^{(1)}, \quad \text{Cov}(\mathbf{x}^{(1)}) = \boldsymbol{\Sigma}_{11}$$

$$E(\mathbf{x}^{(2)}) = \boldsymbol{\mu}^{(2)}, \quad \text{Cov}(\mathbf{x}^{(2)}) = \boldsymbol{\Sigma}_{22}$$

$$\text{Cov}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}_{21}'$$

Será conveniente considerar conjuntamente  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$ , por lo que usando resultados anteriores y la relación anterior, se tiene que

$$\mathbf{X}_{(p+q) \times 1} = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix} = \begin{bmatrix} x_1^{(1)} \\ x_2^{(1)} \\ \vdots \\ x_p^{(1)} \\ x_1^{(2)} \\ x_2^{(2)} \\ \vdots \\ x_q^{(2)} \end{bmatrix}$$

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

$\mathbf{X}$  tiene vector de medias

$$\mu_{(p+q) \times 1} = E[\mathbf{X}] = \begin{bmatrix} \frac{E(\mathbf{x}^{(1)})}{E(\mathbf{x}^{(2)})} \end{bmatrix} = \begin{bmatrix} \mu^{(1)} \\ \mu^{(2)} \end{bmatrix}$$

y matriz de covarianzas

$$\begin{aligned} \Sigma_{(p+q) \times (p+q)} &= E[(\mathbf{X} - \mu)(\mathbf{X} - \mu)'] \\ &= \begin{bmatrix} \frac{E((\mathbf{x}^{(1)} - \mu^{(1)})(\mathbf{x}^{(1)} - \mu^{(1)})')}{E((\mathbf{x}^{(2)} - \mu^{(2)})(\mathbf{x}^{(1)} - \mu^{(1)})')} & \frac{E((\mathbf{x}^{(1)} - \mu^{(1)})(\mathbf{x}^{(2)} - \mu^{(2)})')}{E((\mathbf{x}^{(2)} - \mu^{(2)})(\mathbf{x}^{(2)} - \mu^{(2)})')} \end{bmatrix} \\ &= \begin{bmatrix} \Sigma_{11 \, p \times p} & \Sigma_{12 \, p \times q} \\ \Sigma_{21 \, q \times p} & \Sigma_{22 \, q \times q} \end{bmatrix} \end{aligned}$$

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

- Las covarianzas entre pares de variables de diferentes conjuntos, es decir, una variable de  $\mathbf{x}^{(1)}$  y una variable de  $\mathbf{x}^{(2)}$ , están contenidas en  $\Sigma_{12}$ , o de forma equivalente, en  $\Sigma_{21}$ .
- Es decir, los  $pq$  elementos de  $\Sigma_{12}$  miden la asociación entre los dos conjuntos. Cuando  $p$  y  $q$  son relativamente grandes, la interpretación de los elementos de  $\Sigma_{12}$ , colectivamente resulta inútil y complicada.
- Además, a menudo existen combinaciones lineales de variables que son interesantes y útiles para propósitos predictivos o comparativos.
- La principal tarea del análisis de correlación canónica es resumir las asociaciones entre los conjuntos  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$  en términos de *algunas pocas covarianzas* (o correlaciones) seleccionadas cuidadosamente en lugar de las  $pq$  covarianzas en  $\Sigma_{12}$

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

Las combinaciones lineales proporcionan medidas que resumen de manera sencilla un conjunto de variables. Sean

$$U = \mathbf{a}' \mathbf{x}^{(1)}$$

$$V = \mathbf{b}' \mathbf{x}^{(2)}$$

para algun par de vectores de coeficientes  $\mathbf{a}$  y  $\mathbf{b}$ . Entonces, a partir de lo anterior, obtenemos

$$\text{Var}(U) = \mathbf{a}' \text{Cov}(\mathbf{x}^{(1)}) \mathbf{a} = \mathbf{a}' \boldsymbol{\Sigma}_{11} \mathbf{a}$$

$$\text{Var}(V) = \mathbf{b}' \text{Cov}(\mathbf{x}^{(2)}) \mathbf{b} = \mathbf{b}' \boldsymbol{\Sigma}_{22} \mathbf{b}$$

$$\text{Cov}(U, V) = \mathbf{a}' \text{Cov}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) \mathbf{b} = \mathbf{a}' \boldsymbol{\Sigma}_{12} \mathbf{b}$$

Se buscan los vectores de coeficientes  $\mathbf{a}$  y  $\mathbf{b}$  tal que

$$\text{Corr}(U, V) = \frac{\mathbf{a}' \boldsymbol{\Sigma}_{12} \mathbf{b}}{\sqrt{\mathbf{a}' \boldsymbol{\Sigma}_{11} \mathbf{a}} \sqrt{\mathbf{b}' \boldsymbol{\Sigma}_{22} \mathbf{b}}}$$

sea lo más grande posible



# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

Definimos lo siguiente:

- El primer par de variables canónicas, es el par de combinaciones lineales  $U_1$  y  $V_1$  que tienen varianzas unitarias, que maximizan la correlación anterior.
- El segundo par de variables canónicas es el par de combinaciones lineales  $U_2$  y  $V_2$  que tienen varianzas unitarias, que maximizan la correlación anterior entre todas las elecciones que no están correlacionadas con el primer par de variables canónicas.
- En el  $k$ -ésimo paso, el  $k$ -ésimo par de variables canónicas, es el par de combinaciones lineales  $U_k$  y  $V_k$  que tienen varianzas unitarias, que maximizan la correlación anterior entre todas las elecciones que no están correlacionadas con los anteriores  $k - 1$  pares de variables canónicas. La correlación entre el  $k$ -ésimo par de variables canónicas se denomina  $k$ -ésima correlación canónica.

# Análisis de Correlación Canónica: Variables canónicas y correlaciones Canónicas

**Resultado:** Suponemos que  $p \leq q$  y sean  $\mathbf{x}^{(1)}$  y  $\mathbf{x}^{(2)}$  los vectores aleatorios que tienen  $\text{Cov}(\mathbf{x}^{(1)}) = \mathbf{\Sigma}_{11_{p \times p}}$ ,  $\text{Cov}(\mathbf{x}^{(2)}) = \mathbf{\Sigma}_{22_{q \times q}}$  y  $\text{Cov}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = \mathbf{\Sigma}_{12_{p \times q}}$ , donde  $\mathbf{\Sigma}$  tiene rango completo. Para los coeficientes  $\mathbf{a}_{p \times 1}$  y  $\mathbf{b}_{q \times 1}$ , formamos las combinaciones lineales  $U = \mathbf{a}'\mathbf{x}^{(1)}$  y  $V = \mathbf{b}'\mathbf{x}^{(2)}$ . Entonces

$$\max_{\mathbf{a}, \mathbf{b}} \text{Corr}(U, V) = \rho_1^*$$

se obtiene para las combinaciones lineales (primer par de variables canónicas)

$$U_1 = \mathbf{e}_1' \mathbf{\Sigma}_{11}^{-1/2} \mathbf{x}^{(1)} \quad \text{y} \quad V_1 = \mathbf{f}_1' \mathbf{\Sigma}_{22}^{-1/2} \mathbf{x}^{(2)}$$

# Variables canónicas y correlaciones Canónicas

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

$$U_k = \mathbf{e}_k' \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{x}^{(1)}, \quad V_k = \mathbf{f}_k' \boldsymbol{\Sigma}_{22}^{-1/2} \mathbf{x}^{(2)}$$

maximiza

$$\text{Corr}(U_k, V_k) = \rho_k^*$$

entre aquellas combinaciones lineales no correlacionadas con las anteriores  $1, 2, \dots, k-1$  variables canónicas.

- Aquí  $\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$  son los valores propios de  $\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2}$ , y  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$  de tamaño  $(p \times 1)$ , son los vectores propios asociados
- Las cantidades  $\rho_1^{*2}, \rho_2^{*2}, \dots, \rho_p^{*2}$  son también los  $p$  valores propios mas grandes de la matriz  $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2}$  con correspondientes vectores propios  $(q \times 1)$ ,  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$ .
- Cada  $\mathbf{f}_i$  es proporcional a  $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_i$

Las variables canónicas tienen las siguientes propiedades:

$$\text{Var}(U_k) = \text{Var}(V_k) = 1$$

$$\text{Cov}(U_k, U_l) = \text{Corr}(U_k, U_l) = 0, \quad k \neq l$$

$$\text{Cov}(V_k, V_l) = \text{Corr}(V_k, V_l) = 0, \quad k \neq l$$

$$\text{Cov}(U_k, V_l) = \text{Corr}(U_k, V_l) = 0, \quad k \neq l$$

para  $k, l = 1, 2, \dots, p$

# Variables canónicas y correlaciones Canónicas

Si las variables originales se estandarizan a

$$\mathbf{Z}^{(1)} = [Z_1^{(1)}, Z_2^{(1)}, \dots, Z_p^{(1)}]' \text{ y } \mathbf{Z}^{(2)} = [Z_1^{(2)}, Z_2^{(2)}, \dots, Z_q^{(2)}]',$$

usando los principios anteriores, el  $k$ -ésimo par de variables canónicas son de la forma

$$U_k = \mathbf{a}_k' \mathbf{Z}^{(1)} = \mathbf{e}_k' \boldsymbol{\rho}_{11}^{-1/2} \mathbf{Z}^{(1)}$$

$$V_k = \mathbf{b}_k' \mathbf{Z}^{(2)} = \mathbf{f}_k' \boldsymbol{\rho}_{22}^{-1/2} \mathbf{Z}^{(2)}$$

Donde

- $\text{Cov}(\mathbf{Z}^{(1)}) = \boldsymbol{\rho}_{11}$
- $\text{Cov}(\mathbf{Z}^{(2)}) = \boldsymbol{\rho}_{22}$
- $\text{Cov}(\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}) = \boldsymbol{\rho}_{12} = \boldsymbol{\rho}_{21}'$
- $\mathbf{e}_k$  es el vector propio de  $\boldsymbol{\rho}_{11}^{-1/2} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1/2}$
- $\mathbf{f}_k$  es el vector propio de  $\boldsymbol{\rho}_{22}^{-1/2} \boldsymbol{\rho}_{21} \boldsymbol{\rho}_{11}^{-1} \boldsymbol{\rho}_{12} \boldsymbol{\rho}_{22}^{-1/2}$

Las correlaciones canónicas  $\rho_k^*$  satisfacen

$$\text{Corr}(U_k, V_k) = \rho_k^*, \quad k = 1, 2, \dots, p$$

donde

$$\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$$

son los valores propios distintos de cero de la matriz

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

o equivalentemente los valores propios más grandes de la matriz

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

# Variables canónicas y correlaciones Canónicas

Nótese que

$$\begin{aligned} \mathbf{a}'_k(\mathbf{x}^{(1)} - \boldsymbol{\mu}^{(1)}) &= a_{k1}(x_1^{(1)} - \mu_1^{(1)}) + a_{k2}(x_2^{(1)} - \mu_2^{(1)}) \\ &\quad + \cdots + a_{kp}(x_p^{(1)} - \mu_p^{(1)}) \\ &= a_{k1}\sqrt{\sigma_{11}}\frac{(x_1^{(1)} - \mu_1^{(1)})}{\sqrt{\sigma_{11}}} + a_{k2}\sqrt{\sigma_{22}}\frac{(x_2^{(1)} - \mu_2^{(1)})}{\sqrt{\sigma_{22}}} \\ &\quad + \cdots + a_{kp}\sqrt{\sigma_{pp}}\frac{(x_p^{(1)} - \mu_p^{(1)})}{\sqrt{\sigma_{pp}}} \end{aligned}$$

donde  $\text{Var}(x_i^{(1)}) = \sigma_{ii}, i = 1, 2, \dots, p$ .

Por lo tanto, si  $\mathbf{a}'_k$  es el vector de coeficientes para la  $k$ -ésima variable canónica  $U_k$ , construido a partir de las variables originales  $x_i^{(1)}$  entonces  $\mathbf{a}'_k \mathbf{V}_{11}^{1/2}$  es el vector de coeficientes para la  $k$ -ésima variable canónica construida a partir de las variables estandarizadas  $Z_i^{(1)} = (x_i^{(1)} - \mu_i^{(1)})/\sqrt{\sigma_{ii}}$ , donde  $\mathbf{V}_{11}^{1/2}$  es la matriz diagonal con  $i$ -ésimo elemento diagonal  $\sqrt{\sigma_{ii}}$ .

# Variables canónicas y correlaciones Canónicas

- Similarmente  $\mathbf{b}'_k \mathbf{V}_{22}^{1/2}$  es el vector de coeficientes para la variable canónica construida a partir del conjunto de variables estandarizadas  $\mathbf{Z}^{(2)}$ .
- En este caso  $\mathbf{V}_{22}^{1/2}$  es la matriz diagonal con  $i$ -ésimo elemento diagonal  $\sqrt{\sigma_{ii}} = \sqrt{\text{Var}(x_i^{(2)})}$ .
- Las correlaciones canónicas *no se modifican* por la estandarización de las variables originales, sin embargo, la elección de los vectores de coeficientes  $\mathbf{a}_k, \mathbf{b}_k$  no será única si  $\rho_k^{*2} = \rho_{k+1}^{*2}$ .



- La relación entre los coeficientes canónicos de las variables estandarizadas y los coeficientes canónicos de las variables originales se deriva de la estructura especial de la matriz:

$$\Sigma_{11}^{-1/2} \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21} \Sigma_{11}^{-1/2} \quad \text{o} \quad \rho_{11}^{-1/2} \rho_{12} \rho_{22}^{-1} \rho_{21} \rho_{11}^{-1/2}$$

- y esto es particular al análisis de correlación canónica. Por ejemplo, en el análisis de componentes principales, si  $\mathbf{a}'_k$  es el vector de coeficientes para el  $k$ -ésimo componente principal obtenido de  $\Sigma$  entonces  $\mathbf{a}'_k(\mathbf{X} - \mu) = a'_k \mathbf{V}^{1/2} \mathbf{Z}$ , pero no podemos inferir que  $\mathbf{a}'_k \mathbf{V}^{1/2}$  es el vector de coeficientes para el  $k$ -ésimo componente principal derivado de  $\rho$ .

## Ejemplo: Obtención de variables canónicas y correlaciones canónicas para variables estandarizadas

Supongamos que  $\mathbf{Z}^{(1)} = [Z_1^{(1)}, Z_2^{(1)}]$  y  $\mathbf{Z}^{(2)} = [Z_1^{(2)}, Z_2^{(2)}]$  son variables estandarizadas y  $\mathbf{Z} = [\mathbf{Z}^{(1)}, \mathbf{Z}^{(2)}]$ , donde

$$\text{Cov}(\mathbf{Z}) = \left[ \begin{array}{cc|cc} \rho_{11} & \rho_{12} & & \\ \rho_{21} & \rho_{22} & & \end{array} \right] = \left[ \begin{array}{cc|cc} 1.0 & 0.4 & 0.5 & 0.6 \\ 0.4 & 1.0 & 0.3 & 0.4 \\ \hline 0.5 & 0.3 & 1.0 & 0.2 \\ 0.6 & 0.4 & 0.2 & 1.0 \end{array} \right]$$

Calcular las variables canónicas y correlaciones canónicas.