

# Análise de Correlação Canônica

Prof. George von Borries

Departamento de Estatística  
Universidade de Brasília

2023



## Observação

*O material seguinte é baseado em*

- *Slides de aula de Júlia M. Pavan Soler, Instituto de Matemática e Estatística, USP. (Autorizado pela autora)*
- *Johnson, R.A. e Wichern, D.W. (2007) Applied Multivariate Analysis, 6th Edition. Capítulo 10.*



# Introdução

- Correlação Canônica tem o objetivo de verifica a relação linear **entre** dois **conjuntos** de variáveis.
- Na ACP e AF consideramos a variação **dentro** de um conjunto de variáveis.
- Assumimos que os dois conjuntos de variáveis são medidos na mesma unidade amostral.
- A ideia é determinar o pares de combinações lineares com maior correlação e não correlacionados com pares anteriormente selecionados. Os pares de combinações lineares são chamados de **variáveis canônicas** e suas correlações de **correlação canônica**.
- Exemplos:
  - 1 Relacionar variáveis de política governamental com variáveis econômicas ao longo do tempo.
  - 2 Relacionar variáveis de desempenho escolar no ensino médio e superior.
  - 3 Relacionar variáveis de satisfação numa profissão com variáveis de características da profissão.



- Seja  $\mathbb{R}^{p+q} = \mathbb{R}^p \times \mathbb{R}^q$  tal que  $\mathbf{Y}_{p \times 1} = \{y_1, \dots, y_p\}$  e  $\mathbf{X}_{q \times 1} = \{x_1, \dots, x_q\}$  são dois grupos de variáveis aleatórias, i.e.,

$$\mathbf{M} = \begin{pmatrix} \mathbf{Y} \\ \mathbf{X} \end{pmatrix} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_p \\ x_1 \\ x_2 \\ \vdots \\ x_q \end{pmatrix}$$

- Seguindo a notação de Johnson e Wichern (2007), assumimos  $p \leq q$ .



- Sejam ainda,

$$\boldsymbol{\mu} = E(\mathbf{M}) = \begin{pmatrix} \frac{E(\mathbf{Y})}{E(\mathbf{X})} \end{pmatrix} = \begin{pmatrix} \frac{\mu_y}{\mu_x} \end{pmatrix}$$

e

$$\begin{aligned} \boldsymbol{\Sigma} &= E[(\mathbf{M} - \boldsymbol{\mu})(\mathbf{M} - \boldsymbol{\mu})^T] \\ &= \begin{pmatrix} \frac{E[(\mathbf{Y} - \mu_y)(\mathbf{Y} - \mu_y)^T]}{E[(\mathbf{X} - \mu_x)(\mathbf{Y} - \mu_y)^T]} & \frac{E[(\mathbf{Y} - \mu_y)(\mathbf{X} - \mu_x)^T]}{E[(\mathbf{X} - \mu_x)(\mathbf{X} - \mu_x)^T]} \end{pmatrix} \\ &= \begin{pmatrix} \boldsymbol{\Sigma}_{yy} & \boldsymbol{\Sigma}_{yx} \\ \boldsymbol{\Sigma}_{xy} & \boldsymbol{\Sigma}_{xx} \end{pmatrix} \end{aligned}$$

- $\boldsymbol{\Sigma}_{yx} = \boldsymbol{\Sigma}_{xy}^T$  e mede a associação entre os dois conjuntos de variáveis.
- As covariâncias em  $\boldsymbol{\Sigma}_{yx}$  são medidas em pares, i.e., temos  $p \times q$  covariâncias.



## Combinações Lineares

- A correlação canônica busca obter combinações lineares das variáveis que resumem “o melhor possível” a associação entre as variáveis.
- Sejam então,  $U = \mathbf{a}^T \mathbf{Y}$  e  $V = \mathbf{b}^T \mathbf{X}$ . Então,

$$\begin{aligned}\text{Var}(U) &= \mathbf{a}^T \text{Cov}(\mathbf{Y}) \mathbf{a} = \mathbf{a}^T \boldsymbol{\Sigma}_{yy} \mathbf{a} \\ \text{Var}(V) &= \mathbf{b}^T \text{Cov}(\mathbf{X}) \mathbf{b} = \mathbf{b}^T \boldsymbol{\Sigma}_{xx} \mathbf{b} \\ \text{Cov}(U, V) &= \mathbf{a}^T \text{Cov}(\mathbf{Y}, \mathbf{X}) \mathbf{b} = \mathbf{a}^T \boldsymbol{\Sigma}_{yx} \mathbf{b}\end{aligned}$$

- Procuramos coeficientes  $\mathbf{a}$  e  $\mathbf{b}$  tais que

$$\text{Cor}(U, V) = \frac{\mathbf{a}^T \boldsymbol{\Sigma}_{yx} \mathbf{b}}{\sqrt{\mathbf{a}^T \boldsymbol{\Sigma}_{yy} \mathbf{a}} \sqrt{\mathbf{b}^T \boldsymbol{\Sigma}_{xx} \mathbf{b}}}$$

seja o maior possível.

- O primeiro par de variáveis canônicas  $U_1, V_1$ , tem variância unitária e maximiza a correlação.



- O segundo par de variáveis canônicas  $U_2, V_2$ , tem variância unitária e maximiza a correlação entre todas as possíveis escolhas que são não correlacionadas com o primeiro par de variáveis canônicas.
- O processo continua até obtermos  $p$  correlações de variáveis canônicas.

**Resultado 1:** Suponha  $\Sigma_{(p+q) \times (p+q)}$  como descrito anteriormente e de rank completo. Seja ainda  $p \leq q$ . Então

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

obtido pelo primeiro par de variáveis canônicas

$$\mathbf{U}_1 = \mathbf{a}_1^T \mathbf{Y} = \mathbf{e}_1^T \Sigma_{yy}^{-1/2} \mathbf{Y} \quad \text{e} \quad \mathbf{V}_1 = \mathbf{b}_1^T \mathbf{X} = \mathbf{f}_1^T \Sigma_{xx}^{-1/2} \mathbf{X}$$

e  $\rho_k$  obtido pelo primeiro par de variáveis canônicas

$$\mathbf{U}_k = \mathbf{a}_k^T \mathbf{Y} = \mathbf{e}_k^T \Sigma_{yy}^{-1/2} \mathbf{Y} \quad \text{e} \quad \mathbf{V}_k = \mathbf{b}_k^T \mathbf{X} = \mathbf{f}_k^T \Sigma_{xx}^{-1/2} \mathbf{X}$$

entre as combinações lineares não correlacionadas com os pares  $1, 2, \dots, k-1$  de variáveis canônicas.



## Resultado 1: (continuação)

As correlações  $\rho_1^2$  e  $\mathbf{e}_1$  são o maior autovalor e o respectivo autovetor de

$$\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1/2}$$

e as correlações  $\rho_1^2$  e  $\mathbf{f}_1$  são o maior autovalor e o respectivo autovetor de

$$\Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \Sigma_{xx}^{-1/2}.$$





Generalizando,

- $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$  são autovalores de

$$\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1/2}$$

- $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$  são os autovetores correspondentes ( $p \times 1$ ).
- $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$  são também autovalores de

$$\Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1} \Sigma_{yx} \Sigma_{xx}^{-1/2}.$$

- $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$  são os autovetores correspondentes ( $q \times 1$ ).
- $\mathbf{f}_i$  é proporcional a

$$\Sigma_{xx}^{-1/2} \Sigma_{xy} \Sigma_{yy}^{-1} \mathbf{e}_i$$



Se as variáveis são padronizadas,

$$\text{Cor}(U_k^*, V_k^*) = \frac{\mathbf{a}_k^{*T} \boldsymbol{\rho}_{yx} \mathbf{b}_k^*}{\sqrt{\mathbf{a}_k^{*T} \boldsymbol{\rho}_{yy} \mathbf{a}_k^*} \sqrt{\mathbf{b}_k^{*T} \boldsymbol{\rho}_{xx} \mathbf{b}_k^*}} = \rho_k$$

- **Resultado 2:** As correlações canônicas são invariantes a padronização dos dados.
- $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$  são autovalores de

$$\boldsymbol{\rho}_{yy}^{-1/2} \boldsymbol{\rho}_{yx} \boldsymbol{\rho}_{xx}^{-1} \boldsymbol{\rho}_{xy} \boldsymbol{\rho}_{yy}^{-1/2}$$

- $\mathbf{e}_1^*, \mathbf{e}_2^*, \dots, \mathbf{e}_p^*$  são os autovetores correspondentes ( $p \times 1$ ).
- $\rho_1^2 \geq \rho_2^2 \geq \dots \geq \rho_p^2$  são também autovalores de

$$\boldsymbol{\rho}_{xx}^{-1/2} \boldsymbol{\rho}_{xy} \boldsymbol{\rho}_{yy}^{-1} \boldsymbol{\rho}_{yx} \boldsymbol{\rho}_{xx}^{-1/2}.$$

- $\mathbf{f}_1^*, \mathbf{f}_2^*, \dots, \mathbf{f}_p^*$  são os autovetores correspondentes ( $q \times 1$ ).



- Os coeficientes canônicos, em geral, não tem interpretação física.
- Os coeficientes canônicos das variáveis originais ( $U, V$ ) tem unidades de medida proporcionais àquelas das variáveis originais nos dois grupos.
- Os coeficientes canônicos das variáveis padronizadas ( $U^*, V^*$ ) não tem unidades de medida e devem ser interpretados em função das variáveis padronizadas (média zero e variância unitária).
- Os coeficientes canônicos das variáveis padronizadas podem ser obtidos diretamente dos coeficientes das variáveis originais.
- O coeficiente canônico das variáveis originais e das variáveis padronizadas é o mesmo, i.e., o coeficiente de correlação canônico é invariante sob padronização das variáveis.



## Interpretação Geométrica

$$U_1 = \mathbf{a}_1^T \mathbf{Y} = \mathbf{e}_1^T \Sigma_{yy}^{-1/2} \mathbf{Y} = \mathbf{e}_1^T \mathbf{P}_1 \underbrace{\mathbf{D}^{-1/2} \mathbf{P}_1^T \mathbf{Y}}_{\text{Fator Comum de } \mathbf{Y}} \quad \begin{matrix} \text{CP de } \mathbf{Y} \\ \mathbf{P}_1^T \mathbf{Y} \end{matrix}$$

- A variável canônica  $U_1$  resulta de uma rotação ortogonal (via  $\mathbf{P}_1$  e determinada por  $\Sigma_{yy}$ ) do CP padronizado (Fator Comum de  $\mathbf{Y}$ ) seguida por outra rotação ortogonal (via  $\mathbf{e}_1$  e determinada por  $\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1/2}$ ).
- Se  $\lambda$  é autovalor de  $\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy} \Sigma_{yy}^{-1/2}$  com  $\mathbf{e}$  o correspondente autovetor, então  $\lambda$  é também autovalor de  $\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$  com  $\Sigma_{yy}^{-1/2} \mathbf{e}$  o correspondente autovetor.
- Assim, as variáveis e os coeficientes canônicos podem ser obtidos diretamente da decomposição em valores singulares de  $\Sigma_{yy}^{-1/2} \Sigma_{yx} \Sigma_{xx}^{-1} \Sigma_{xy}$ .



## Exemplos de Johnson e Wichern - Capítulo 10:

J&W-Exemplos-Cap10.R

- Exemplo 10.1 - pág. 543 a 545.
- Exemplo 10.2 - pág. 546 a 547.
- Exemplo 10.3 - pág. 549.
- Exemplo 10.4 - pág. 552 e 553.
- Exemplo 10.5 - pág. 553 a 555.

