



## Unidade VII – Análise de Correlação Canônica

### 7.1. Introdução

A análise de correlação canônica (ACC) é a técnica multivariada que estuda o relacionamento (linear) entre dois conjuntos de variáveis. Em muitos estudos são avaliados dois tipos de variáveis em cada unidade de pesquisa, por exemplo, um conjunto de variáveis de aptidão e um conjunto de variáveis de desempenho, um conjunto de variáveis de personalidade e um conjunto de medidas de habilidade, um conjunto de índices de preços e um conjunto de índices de produção, um conjunto de comportamentos dos alunos e um conjunto de comportamentos do professor, um conjunto de atributos psicológicos e um conjunto de atributos fisiológicos, um conjunto de variáveis ecológicas e um conjunto de variáveis ambientais, um conjunto de desempenho acadêmico e um conjunto de medidas de sucesso no trabalho, e um conjunto de variáveis de personalidade de alunos ao ingressarem na faculdade e as mesmas variáveis nos mesmos indivíduos quando idosos.

Na ACC procura-se identificar e quantificar a associação entre esses dois conjuntos de variáveis por meio do desenvolvimento de uma combinação linear das variáveis em cada um dos grupos, de modo que a correlação entre essas duas combinações seja maximizada. Essas combinações lineares são as variáveis canônicas e suas associações são denominadas de correlações canônicas.

A ACC objetiva encontrar combinações lineares que expressem bem as correlações entre os dois conjuntos de variáveis e obter a simplificação dos dados, ao descrever os dados em um número muito menor variáveis, apontando quais as variáveis originais são as mais importantes.

### 7.2. Modelo Teórico

Seja  $\mathbf{X}$  um vetor de dimensão  $(p+q) \times 1$  que possui vetor de médias  $\boldsymbol{\mu}$  e matriz de covariância  $\boldsymbol{\Sigma}$ , que é positiva definida. Consideram-se dois conjuntos de variáveis em  $\mathbf{X}$ . O 1º conjunto,  $\mathbf{X}^{(1)}$ , com  $p$  variáveis, e o 2º,  $\mathbf{X}^{(2)}$ , com  $q$  variáveis. Sem perda de generalidade assume-se que  $p \leq q$ . A matriz  $\mathbf{X}$ , o vetor de médias  $\boldsymbol{\mu}$  e a matriz de covariância  $\boldsymbol{\Sigma}$  podem ser particionados da seguinte maneira:

$$\mathbf{X}'_{1 \times (p+q)} = \begin{bmatrix} \mathbf{X}^{(1)} & \mathbf{X}^{(2)} \end{bmatrix}_{(1 \times p) \quad (1 \times q)}, \quad \boldsymbol{\mu}'_{1 \times (p+q)} = \begin{bmatrix} \boldsymbol{\mu}^{(1)} & \boldsymbol{\mu}^{(2)} \end{bmatrix}_{(1 \times p) \quad (1 \times q)}, \quad \boldsymbol{\Sigma}_{(p+q) \times (p+q)} = E(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})' = \begin{bmatrix} \boldsymbol{\Sigma}_{11} & \vdots & \boldsymbol{\Sigma}_{12} \\ \boldsymbol{\Sigma}_{21} & \vdots & \boldsymbol{\Sigma}_{22} \end{bmatrix}$$

$\begin{matrix} (p \times p) & & (p \times q) \\ \dots & & \dots \\ (q \times p) & & (q \times q) \end{matrix}$

onde

$\boldsymbol{\Sigma}_{11} = Cov(\mathbf{X}^{(1)})$  é a matriz de covariância para o 1º conjunto de dados.

$\boldsymbol{\Sigma}_{22} = Cov(\mathbf{X}^{(2)})$  é a matriz de covariância para o 2º conjunto de dados.

$\boldsymbol{\Sigma}_{12} = \boldsymbol{\Sigma}'_{21} = Cov(\mathbf{X}^{(1)}, \mathbf{X}^{(2)})$  é a matriz de covariância entre os dois conjuntos de dados, ela mede a associação entre os dois conjuntos.

Como o maior interesse está na matriz  $\boldsymbol{\Sigma}_{12}$ , mas ela pode ser de grande dimensão e, portanto, difícil de analisar, a ideia da técnica é estudar algumas poucas combinações lineares de variáveis pertencentes a  $\mathbf{X}^{(1)}$  e  $\mathbf{X}^{(2)}$ , ao invés de usar a matriz de covariância  $\boldsymbol{\Sigma}_{12}$ .

São definidas como  $U$  e  $V$  as combinações lineares das variáveis de  $\mathbf{X}^{(1)}$  e de  $\mathbf{X}^{(2)}$ :

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

Sendo  $\mathbf{a}$  e  $\mathbf{b}$  vetores não nulos dos coeficientes dessas combinações lineares.

As variâncias e a covariância entre  $U$  e  $V$  são definidas como:

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

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

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

E a correlação entre  $U$  e  $V$  é definida como:

$$Cor(U, V) = \rho_{U,V} = \frac{\mathbf{a}'\boldsymbol{\Sigma}_{12}\mathbf{b}}{\sqrt{\mathbf{a}'\boldsymbol{\Sigma}_{11}\mathbf{a}}\sqrt{\mathbf{b}'\boldsymbol{\Sigma}_{22}\mathbf{b}}} \quad (6.1)$$

Hotelling (1935) propôs estimar os vetores  $\mathbf{a}$  e  $\mathbf{b}$  maximizando  $Cor(U, V)$  sujeita à restrição:  $\mathbf{a}'\boldsymbol{\Sigma}_{11}\mathbf{a} = \mathbf{b}'\boldsymbol{\Sigma}_{22}\mathbf{b} = 1$ . As variáveis  $U$  e  $V$  são chamadas variáveis canônicas e a correlação entre elas, correlação canônica.

Assumindo  $p < q$ , teremos  $p$  correlações canônicas que ao serem maximizadas irão gerar  $p$  variáveis canônicas em  $\mathbf{X}^{(1)}$  chamadas de  $U_k$  e  $p$  variáveis canônicas em  $\mathbf{X}^{(2)}$  chamadas de  $V_k$ . Então, o primeiro par de variáveis canônicas (primeira função canônica), é o par de combinações lineares  $U_1$  e  $V_1$ , com variâncias unitárias, que maximiza a correlação (6.1). A segunda função canônica, é o par de combinações lineares  $U_2$  e  $V_2$ , com variâncias unitárias, que maximiza a correlação (6.1) dentre todas as possíveis funções canônicas não correlacionadas com a primeira. E assim segue até o par  $p$ , pois são encontradas tantas funções canônicas quantos forem o menor número entre  $p$  e  $q$ . A  $p$ -ésima função canônica deve ser não correlacionada com as  $(p-1)$  funções canônicas anteriores.

Cada par de variáveis canônicas,  $U_k$  e  $V_k$ , que apresentam variância unitária, é definido pelos seus respectivos vetores,  $(\mathbf{a}_k$  e  $\mathbf{b}_k)$ , que maximizam  $Cor(U_k, V_k) = \rho_k^*$ .

O resultado da maximização são as combinações lineares do  $k$ -ésimo par de variáveis canônicas:

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

$$\text{tal que } \text{Max } Cor(U_k, V_k) = \rho_k^* = \sqrt{\lambda_k} \Rightarrow \lambda_k = (\mathbf{a}_k' \boldsymbol{\Sigma}_{12} \mathbf{b}_k)^2$$

$$\text{Em que } \lambda_k \text{ satisfaz: } \begin{cases} \left( (\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} - \lambda_k \mathbf{I}) \mathbf{e}_k = 0 \right. \\ \left. \left( \boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2} - \lambda_k \mathbf{I} \right) \mathbf{f}_k = 0 \right. \end{cases}$$

Então  $\rho_1^{*2} \geq \rho_2^{*2} \geq \dots \geq \rho_p^{*2}$  são os  $p$  maiores autovalores da matriz  $\boldsymbol{\Sigma}_{11}^{-1/2} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2}$  e  $\mathbf{e}_1, \mathbf{e}_2, \dots, \mathbf{e}_p$  são os autovetores associados. As quantidades  $\rho_1^{*2}, \rho_2^{*2}, \dots, \rho_p^{*2}$  são também os  $p$  maiores autovalores da matriz  $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1} \boldsymbol{\Sigma}_{12} \boldsymbol{\Sigma}_{22}^{-1/2}$  com os correspondentes autovetores  $\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_p$ . Cada  $\mathbf{f}_i$  é proporcional a  $\boldsymbol{\Sigma}_{22}^{-1/2} \boldsymbol{\Sigma}_{21} \boldsymbol{\Sigma}_{11}^{-1/2} \mathbf{e}_i$ . As variáveis canônicas têm as seguintes propriedades:

$$\begin{aligned} Var(U_k) &= Var(V_k) = 1 \\ Cov(U_k, U_l) &= Cov(V_k, V_l) = Cov(U_k, V_l) = Cov(U_l, V_k) = 0 \quad (k \neq l) \\ Corr(U_k, U_l) &= Corr(V_k, V_l) = Corr(U_k, V_l) = Corr(U_l, V_k) = 0 \quad (k \neq l) \end{aligned}$$

Para o caso amostral, a matriz de covariância particionada  $\boldsymbol{\Sigma}$  deverá ser substituída pela sua estimativa amostral,  $\mathbf{S}$ .

Se as variáveis forem padronizadas não vai alterar a correlação canônica e o sistema para estimar os coeficientes  $\mathbf{a}_k$  e  $\mathbf{b}_k$  terá solução única se a matriz  $\boldsymbol{\Sigma}$  tiver posto completo.

### 7.3. Interpretação das variáveis canônicas

Existem três formas de interpretação das variáveis canônicas. O nível de significância das raízes canônicas, a magnitude das correlações canônicas e as relações entre as variâncias das variáveis canônicas e variáveis originais. Essas relações podem ser de três tipos: os próprios pesos canônicos que estamos chamando de coeficientes canônicos, as cargas canônicas ou estrutura de correlação canônica e as cargas canônicas cruzadas.

As correlações entre as variáveis canônicas e seus respectivos conjuntos originais são chamadas de estrutura canônica e as correlações entre as variáveis canônicas e os conjuntos de variáveis opostas são chamadas de cargas cruzadas. Essas correlações devem ser interpretadas com cuidado porque elas trazem informações univariadas, ou seja, não indicam como a variável original contribui conjuntamente para a análise canônica. Por essa razão, muitos pesquisadores preferem avaliar a contribuição das variáveis originais diretamente pelos coeficientes canônicos calculados a partir das variáveis padronizadas (coeficientes canônicos padronizados).

Considerando as matrizes  $\mathbf{A}_{(p \times p)} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_p]'$  e  $\mathbf{B}_{(q \times q)} = [\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_q]'$ , cujas linhas são os vetores de coeficientes das variáveis canônicas, os vetores de variáveis canônicas são:

$$\mathbf{U}_{(p \times 1)} = \mathbf{A} \mathbf{X}^{(1)} \text{ e } \mathbf{V}_{(q \times 1)} = \mathbf{B} \mathbf{X}^{(2)} \quad (6.2)$$

O interesse maior são as primeiras  $p$  variáveis canônicas em  $\mathbf{V}$ . Então:

$$Cov(\mathbf{U}, \mathbf{X}^{(1)}) = Cov(\mathbf{A} \mathbf{X}^{(1)}, \mathbf{X}^{(1)}) = \mathbf{A} Cov(\mathbf{X}^{(1)}) = \mathbf{A} \boldsymbol{\Sigma}_{11}$$

Como cada  $U_i$  tem variância unitária,

$$Corr(U_i, X_k^{(1)}) = \frac{Cov(U_i, X_k^{(1)})}{\sqrt{Var(U_i)} \sqrt{Var(X_k^{(1)})}} = \frac{Cov(U_i, X_k^{(1)})}{\sqrt{Var(X_k^{(1)})}} = \frac{Cov(U_i, X_k^{(1)})}{\sigma_{kk}^{1/2}} = Cov(U_i, \sigma_{kk}^{-1/2} X_k^{(1)})$$

Para colocar a correlação na forma matricial introduzimos a matriz diagonal  $p \times p$ ,  $\mathbf{V}_{11}^{-1/2}$  que possui como  $k$ -ésimo elemento  $\sigma_{kk}^{-1/2}$ :

$$\rho_{u,x^{(1)}} = \text{Corr}(\mathbf{u}, \mathbf{x}^{(1)}) = \text{Cov}\left(\mathbf{u}, \mathbf{V}_{11}^{-1/2} \mathbf{x}^{(1)}\right) = \text{Cov}\left(\mathbf{A} \mathbf{x}^{(1)}, \mathbf{V}_{11}^{-1/2} \mathbf{x}^{(1)}\right) = \mathbf{A} \text{Cov}(\mathbf{x}^{(1)}, \mathbf{x}^{(1)}) \mathbf{V}_{11}^{-1/2} = \mathbf{A} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2}$$

(p x p)

Cálculos similares podem ser realizados para determinar os outros pares de correlações:

$$\begin{aligned} \rho_{u,x^{(1)}} &= \mathbf{A} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2} & \rho_{v,x^{(2)}} &= \mathbf{B} \boldsymbol{\Sigma}_{22} \mathbf{V}_{22}^{-1/2} \\ & \text{(p x p)} & & \text{(q x q)} \\ \rho_{u,x^{(2)}} &= \mathbf{A} \boldsymbol{\Sigma}_{12} \mathbf{V}_{22}^{-1/2} & \rho_{v,x^{(1)}} &= \mathbf{B} \boldsymbol{\Sigma}_{21} \mathbf{V}_{11}^{-1/2} \\ & \text{(p x q)} & & \text{(q x p)} \end{aligned}$$

onde  $\mathbf{V}_{22}^{-1/2}$  é a matriz diagonal ( $q \times q$ ) que tem  $\text{Var}(X_i^{(2)})$  como o  $i$ -ésimo elemento da diagonal.

Se as variáveis forem padronizadas tem-se:

$$\begin{aligned} \rho_{u,z^{(1)}} &= \mathbf{A}_z \boldsymbol{\rho}_{11} & \rho_{v,z^{(2)}} &= \mathbf{B}_z \boldsymbol{\rho}_{22} \\ & \text{(p x p)} & & \text{(q x q)} \\ \rho_{u,z} &= \mathbf{A}_z \boldsymbol{\rho}_{12} & \rho_{v,z} &= \mathbf{B}_z \boldsymbol{\rho}_{21} \\ & \text{(p x q)} & & \text{(q x p)} \end{aligned}$$

onde  $\mathbf{A}_z$  e  $\mathbf{B}_z$  são as matrizes cujas linhas contém os coeficientes canônicos para os conjuntos de variáveis padronizadas  $\mathbf{Z}^{(1)}$  e  $\mathbf{Z}^{(2)}$  respectivamente.

As correlações podem então ser classificadas como:

a) Estrutura canônica:

$$\begin{aligned} \rho_{u,x^{(1)}} &= \mathbf{A} \boldsymbol{\Sigma}_{11} \mathbf{V}_{11}^{-1/2} & \text{e} & & \rho_{v,x^{(2)}} &= \mathbf{B} \boldsymbol{\Sigma}_{22} \mathbf{V}_{22}^{-1/2} \\ & \text{(p x p)} & & & \text{(q x q)} \\ & & \text{ou} & & \\ \rho_{u,z^{(1)}} &= \mathbf{A}_z \boldsymbol{\rho}_{11} & \text{e} & & \rho_{v,x^{(2)}} &= \mathbf{B}_z \boldsymbol{\rho}_{22} \\ & \text{(p x p)} & & & \text{(q x q)} \end{aligned}$$

b) Cargas canônicas cruzadas:

$$\begin{aligned} \rho_{u,x^{(2)}} &= \mathbf{A} \boldsymbol{\Sigma}_{12} \mathbf{V}_{22}^{-1/2} & \text{e} & & \rho_{v,x^{(1)}} &= \mathbf{B} \boldsymbol{\Sigma}_{21} \mathbf{V}_{11}^{-1/2} \\ & \text{(p x q)} & & & \text{(q x p)} \\ & & \text{ou} & & \\ \rho_{u,x^{(2)}} &= \mathbf{A}_z \boldsymbol{\rho}_{12} & \text{e} & & \rho_{v,x^{(1)}} &= \mathbf{B}_z \boldsymbol{\rho}_{21} \\ & \text{(p x q)} & & & \text{(q x p)} \end{aligned}$$

#### 7.4. Proporção da Variância explicada e Índice de redundância

A correlação canônica elevada ao quadrado representa uma estimativa da variância conjunta (ou compartilhada) entre as variáveis canônicas. É uma medida que pode ser mal interpretada uma vez que mede a variância compartilhada entre as variáveis canônicas e não entre as variáveis originais.

A **proporção da variância explicada** é a proporção da variância de um conjunto de variáveis que é explicada pelas respectivas variáveis canônicas. Dessa forma, mesmo que as correlações canônicas sejam fortes podem não ter sido extraídas quantidades significativas da variância das variáveis originais.

O **índice de redundância** foi proposto para facilitar as interpretações. Ele é equivalente ao coeficiente de determinação da análise de regressão. É a média simples dos coeficientes de correlação múltiplo de um conjunto de variáveis com cada uma das variáveis do outro conjunto, que resulta num coeficiente de determinação médio. Essa medida mede a porcentagem da variância em um conjunto de variáveis que é explicada pelo outro conjunto. Deve-se notar que o valor máximo desse coeficiente não é 100% e sim a variância compartilhada entre os dois conjuntos.

O índice de redundância pode ser calculado para cada variável canônica e depois ser calculado o índice de redundância total, que nada mais é do que a quantidade de variância em um conjunto explicado pelas  $r$  primeiras variáveis canônicas do outro.

Quando as observações são padronizadas, as matrizes de covariâncias amostrais  $\mathbf{S}_{kl}$  são as matrizes de correlação  $\mathbf{R}_{kl}$ . Os vetores de coeficientes canônicos são as linhas das matrizes  $\hat{\mathbf{A}}_z$  e  $\hat{\mathbf{B}}_z$  e as colunas de  $\hat{\mathbf{A}}_z^{-1}$  e  $\hat{\mathbf{B}}_z^{-1}$  são as correlações amostrais entre as variáveis canônicas e as variáveis originais padronizadas que as compõem.

Especificamente:

$$\begin{aligned} \text{Cov}(\mathbf{z}^{(1)}, \hat{\mathbf{U}}) &= \text{Cov}(\hat{\mathbf{A}}_z^{-1} \hat{\mathbf{U}}, \hat{\mathbf{U}}) = \hat{\mathbf{A}}_z^{-1} \mathbf{I} = \hat{\mathbf{A}}_z^{-1} \\ \text{Cov}(\mathbf{z}^{(1)}, \hat{\mathbf{V}}) &= \text{Cov}(\hat{\mathbf{B}}_z^{-1} \hat{\mathbf{V}}, \hat{\mathbf{V}}) = \hat{\mathbf{B}}_z^{-1} \mathbf{I} = \hat{\mathbf{B}}_z^{-1} \end{aligned}$$

Então,

$$\hat{\mathbf{A}}_z^{-1} = [\hat{\mathbf{a}}_z^{(1)}, \hat{\mathbf{a}}_z^{(2)}, \dots, \hat{\mathbf{a}}_z^{(p)}] = \begin{bmatrix} r_{\hat{U}_1, z_1^{(1)}} & r_{\hat{U}_2, z_1^{(1)}} & \cdots & r_{\hat{U}_p, z_1^{(1)}} \\ r_{\hat{U}_1, z_2^{(1)}} & r_{\hat{U}_2, z_2^{(1)}} & \cdots & r_{\hat{U}_p, z_2^{(1)}} \\ \vdots & \vdots & \ddots & \vdots \\ r_{\hat{U}_1, z_p^{(1)}} & r_{\hat{U}_2, z_p^{(1)}} & \cdots & r_{\hat{U}_p, z_p^{(1)}} \end{bmatrix}$$

e

$$\hat{\mathbf{B}}_z^{-1} = [\hat{\mathbf{b}}_z^{(1)}, \hat{\mathbf{b}}_z^{(2)}, \dots, \hat{\mathbf{b}}_z^{(p)}] = \begin{bmatrix} r_{\hat{V}_1, z_1^{(2)}} & r_{\hat{V}_2, z_1^{(2)}} & \cdots & r_{\hat{V}_p, z_1^{(2)}} \\ r_{\hat{V}_1, z_2^{(2)}} & r_{\hat{V}_2, z_2^{(2)}} & \cdots & r_{\hat{V}_p, z_2^{(2)}} \\ \vdots & \vdots & \ddots & \vdots \\ r_{\hat{V}_1, z_p^{(2)}} & r_{\hat{V}_2, z_p^{(2)}} & \cdots & r_{\hat{V}_p, z_p^{(2)}} \end{bmatrix}$$

onde  $r_{\hat{U}_i, z_i^{(1)}}$  e  $r_{\hat{V}_i, z_i^{(2)}}$  são os coeficientes de correlação amostrais entre as variáveis originais e as variáveis canônicas, ou seja, os termos que compõem a matriz de estrutura canônica.

Para observações padronizadas, tem-se:

Variância total amostral (padronizada) no primeiro conjunto de variáveis =

$$= \text{tr}(\mathbf{R}_{11}) = \text{tr} [\hat{\mathbf{a}}_z^{(1)} \hat{\mathbf{a}}_z^{(1)'} + \hat{\mathbf{a}}_z^{(2)} \hat{\mathbf{a}}_z^{(2)'} + \cdots + \hat{\mathbf{a}}_z^{(p)} \hat{\mathbf{a}}_z^{(p)'}] = p$$

Variância total amostral (padronizada) no segundo conjunto de variáveis =

$$= \text{tr}(\mathbf{R}_{22}) = \text{tr} [\hat{\mathbf{b}}_z^{(1)} \hat{\mathbf{b}}_z^{(1)'} + \hat{\mathbf{b}}_z^{(2)} \hat{\mathbf{b}}_z^{(2)'} + \cdots + \hat{\mathbf{b}}_z^{(p)} \hat{\mathbf{b}}_z^{(p)'}] = q$$

Como as correlações nas primeiras  $r < p$  colunas de  $\hat{\mathbf{A}}_z^{-1}$  e  $\hat{\mathbf{B}}_z^{-1}$  envolvem só as variáveis canônicas amostrais,  $[\hat{U}_1, \hat{U}_2, \dots, \hat{U}_r]$  e  $[\hat{V}_1, \hat{V}_2, \dots, \hat{V}_r]$  respectivamente, definimos a contribuição das primeiras  $r$  variáveis canônicas à variância total amostral (padronizada) como:

$$\begin{aligned} \text{tr} [\hat{\mathbf{a}}_z^{(1)} \hat{\mathbf{a}}_z^{(1)'} + \hat{\mathbf{a}}_z^{(2)} \hat{\mathbf{a}}_z^{(2)'} + \cdots + \hat{\mathbf{a}}_z^{(r)} \hat{\mathbf{a}}_z^{(r)'}] &= \sum_{i=1}^r \sum_{k=1}^p r_{\hat{U}_i, z_k^{(1)}}^2 \\ \text{tr} [\hat{\mathbf{b}}_z^{(1)} \hat{\mathbf{b}}_z^{(1)'} + \hat{\mathbf{b}}_z^{(2)} \hat{\mathbf{b}}_z^{(2)'} + \cdots + \hat{\mathbf{b}}_z^{(p)} \hat{\mathbf{b}}_z^{(p)'}] &= \sum_{i=1}^r \sum_{k=1}^p r_{\hat{V}_i, z_k^{(2)}}^2 \end{aligned}$$

Logo, a proporção da variância total amostral (padronizada) explicada pelas  $r$  primeiras variáveis canônicas  $(\hat{U}_1, \hat{U}_2, \dots, \hat{U}_r)$  para o 1º conjunto é dada por:

$$\mathbf{R}_{z^{(1)}|\hat{U}_1, \hat{U}_2, \dots, \hat{U}_r}^2 = \frac{\text{tr} [\hat{\mathbf{a}}_z^{(1)} \hat{\mathbf{a}}_z^{(1)'} + \hat{\mathbf{a}}_z^{(2)} \hat{\mathbf{a}}_z^{(2)'} + \cdots + \hat{\mathbf{a}}_z^{(r)} \hat{\mathbf{a}}_z^{(r)'}]}{\text{tr}(\mathbf{R}_{11})} = \frac{\sum_{i=1}^r \sum_{k=1}^p r_{\hat{U}_i, z_k^{(1)}}^2}{p}$$

que representa a soma de cada valor da matriz de estrutura canônica elevado ao quadrado dividida por  $p$  (que é o número de variáveis no primeiro conjunto).

De forma análoga, para o segundo conjunto de variáveis, a proporção da variância total amostral padronizada explicada por  $\hat{V}_1, \hat{V}_2, \dots, \hat{V}_r$  é dada por:

$$\mathbf{R}_{z^{(2)}|\hat{V}_1, \hat{V}_2, \dots, \hat{V}_r}^2 = \frac{\text{tr} [\hat{\mathbf{b}}_z^{(1)} \hat{\mathbf{b}}_z^{(1)'} + \hat{\mathbf{b}}_z^{(2)} \hat{\mathbf{b}}_z^{(2)'} + \cdots + \hat{\mathbf{b}}_z^{(p)} \hat{\mathbf{b}}_z^{(p)'}]}{\text{tr}(\mathbf{R}_{22})} = \frac{\sum_{i=1}^r \sum_{k=1}^p r_{\hat{V}_i, z_k^{(2)}}^2}{q}$$

que representa a soma de cada valor da matriz de estrutura canônica elevado ao quadrado dividida por  $q$  (que é o número de variáveis no segundo conjunto).

O índice de redundância das combinações lineares estimadas é determinado como:

$$IR_{U_i} = \hat{\rho}_i^2 \times \frac{\sum_{i=1}^p r_{\hat{\theta}_{i,z_i}^{(1)}}^2}{p} = \lambda_i \times \frac{\sum_{i=1}^p r_{\hat{\theta}_{i,z_i}^{(1)}}^2}{p}$$

ou seja, é o autovalor vezes a proporção da variância explicada pelo respectivo autovalor.

Da mesma forma para o segundo conjunto:

$$IR_{V_i} = \hat{\rho}_i^2 \times \frac{\sum_{i=1}^p r_{\hat{\nu}_{i,z_i}^{(2)}}^2}{q} = \lambda_i \times \frac{\sum_{i=1}^p r_{\hat{\nu}_{i,z_i}^{(2)}}^2}{q}$$

O índice de redundância total nada mais é do que a soma dos índices de redundância de cada combinação linear do mesmo conjunto e dados, ou seja:

$$IRT_U = \sum_{i=1}^p \hat{\rho}_i^2 \times \frac{\sum_{i=1}^p r_{\hat{\theta}_{i,z_i}^{(1)}}^2}{p} = \sum_{i=1}^p \lambda_i \times \frac{\sum_{i=1}^p r_{\hat{\theta}_{i,z_i}^{(1)}}^2}{p}$$

$$IRT_V = \sum_{i=1}^p \hat{\rho}_i^2 \times \frac{\sum_{i=1}^p r_{\hat{\nu}_{i,z_i}^{(2)}}^2}{q} = \sum_{i=1}^p \lambda_i \times \frac{\sum_{i=1}^p r_{\hat{\nu}_{i,z_i}^{(2)}}^2}{q}$$

## 7.5. Testes de significância

Quando se trabalha com grandes amostras há interesse em realizar inferências dos resultados da ACC. Como qualquer teste estatístico necessita-se saber qual a significância de cada correlação canônica. Existem testes de significância global e o mais utilizado é o teste de Rao. Para testar separadamente cada função canônica existem alguns testes. São eles: Lambda de Wilk (Wilk's Lambda), traço de Hotteling-Lawley (Hotteling's trace), Traço de Pillai (Pillai's trace) e maior raiz característica de Roy (Roy's gnc).

Quando  $\Sigma_{12} = 0$  ( $\rho_1^* = \rho_2^* = \dots = \rho_p^* = 0$ ),  $\mathbf{a}'\mathbf{X}^{(1)}$  e  $\mathbf{b}'\mathbf{X}^{(2)}$  têm covariância  $\mathbf{a}'\Sigma_{12}\mathbf{b} = 0$  para todos os vetores  $\mathbf{a}$  e  $\mathbf{b}$ . Consequentemente, todas as correlações canônicas serão zero, e não existe porque propor uma análise de correlação canônica. Pode-se então testar

$$H_0: \Sigma_{12} = 0 \quad \text{versus} \quad H_1: \Sigma_{12} \neq 0$$

utilizando o teste de razão de verossimilhança:

$$T = -2\ln(\Lambda) = n\ln\left(\frac{|\mathbf{S}_{11}||\mathbf{S}_{22}|}{|\mathbf{S}|}\right) = -n\ln\prod_{i=1}^p (1 - \lambda_i) \sim \chi_{pq}^2$$

Rejeita-se a  $H_0$  se o  $T \geq T_c$  de uma distribuição qui-quadrado com  $pq$  graus de liberdade, dado nível de significância ( $\alpha$ ).

Se o teste global rejeitar  $H_0$ , pode ser de interesse avaliar se as  $k$  primeiras correlações são significativas e, então, se as variáveis canônicas seriam importantes para a caracterização dos dois conjuntos de dados. As hipóteses a serem testadas são:

$$H_0^k = (\rho_1^*)^2 \neq 0, (\rho_2^*)^2 \neq 0, \dots, (\rho_k^*)^2 \neq 0, (\rho_{k+1}^*)^2 \neq 0, \dots, (\rho_p^*)^2 \neq 0$$

$$H_1^k = (\rho_i^*)^2 = 0, \text{ para algum } i \geq k + 1$$

Que podem ser testadas utilizando a estatística:

$$T = -\left(n - 1 - \frac{1}{2}(p + q + 1)\right) \ln \prod_{i=k+1}^p (1 - \rho_i^2) \sim \chi_{(p-k+1)(q-k+1)}^2$$