

Análise Multivariada

Análise de Correspondências

Prof. George von Borries

Departamento de Estatística
Universidade de Brasília

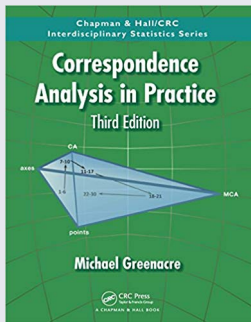
2023



Notas de aula baseadas em

Greenacre, M. (2008) La práctica del análisis de correspondências. Fundación BBVA. Disponível em <http://www.car-me-n.org> **Inzeman, A.J. (2008)** *Modern Multivariate Statistical Techniques*, Springer. **Pavan, J. (2017)** Análise de Correspondência. Slides de notas de aula. IME-USP.

Kassambara, A. Practical Guide to Principal Component Methods in R (Multivariate Analysis II). sthda.com



Definição

Análise de Correspondências (CA): *técnica exploratória que descreve graficamente os dados dispostos em tabelas de contingência para representar o padrão de associação entre linhas e colunas.*

- **Tabela de Contingência** é uma tabela que apresenta o número de observações para múltiplas variáveis categóricas, dispostas em linhas e colunas desta tabela.

Tabela: Relação entre a Cor dos Olhos e Cor do Cabelo entre Crianças Escocesas*.

Olhos	Cabelos				
	CA	CB	CC	CD	CE
OA	621	204	1894	44	70
OB	1002	290	3148	92	139
OC	2827	655	5032	314	150
OD	909	235	4176	226	106

* Exemplo de Izenman (2008). A \rightarrow E indica claro \rightarrow escuro.

- **CA Simples:** estuda tabelas de contingência de dupla entrada, i.e., com duas variáveis dispostas em várias categorias.
- **CA Múltipla:** estuda tabelas de contingência de ordem ≥ 3 .



Perfis e Espaço de Perfis

Tabela: Tabela de Contingência

Var. Linha	Var. Coluna					Total
	1	...	j	...	c	
1	n_{11}		n_{1j}		n_{1c}	$n_{1.}$
\vdots	\vdots	\vdots
i	n_{i1}		n_{ij}		n_{ic}	$n_{i.}$
\vdots	\vdots	\vdots
r	n_{r1}		n_{rj}		n_{rc}	$n_{r.}$
Total	$n_{.1}$		$n_{.j}$		$n_{.c}$	n

Os perfis de linha e coluna serão, respectivamente,

$$p_{ij}^r = \frac{n_{ij}}{n_{i.}} \quad i = 1, \dots, r \quad \text{e} \quad p_{ij}^c = \frac{n_{ij}}{n_{.j}} \quad j = 1, \dots, c.$$



Tabela: Perfis Linha

Var. Linha	Var. Coluna			Total
	1	...	c	
1	$p_{11}^1 = n_{11}/n_{1.}$...	$p_{1c}^1 = n_{1c}/n_{1.}$	1
\vdots	\vdots	\vdots
r	$p_{r1}^r = n_{r1}/n_{r.}$...	$p_{rc}^r = n_{rc}/n_{r.}$	1

Tabela: Perfis Coluna

Var. Linha	Var. Coluna		
	1	...	c
1	$p_{11}^1 = n_{11}/n_{.1}$...	$p_{1c}^c = n_{1c}/n_{.c}$
\vdots	\vdots
r	$p_{r1}^1 = n_{r1}/n_{.1}$...	$p_{rc}^J = n_{rc}/n_{.c}$
Total	1	...	1



Massas e Centróides

- $n = \sum_i n_{i.} = \sum_j n_{.j} = \sum_i \sum_j n_{ij}$ = total de observações.
- Massas de Linhas e Colunas (perfis linha/coluna médio):

$$\frac{n_{i.}}{n} \quad i = 1, \dots, r \quad \frac{n_{.j}}{n} \quad j = 1, \dots, c.$$

Indica o peso de cada perfil (linha/coluna).

- Centróides:

$$\bar{p}^r = \left(\frac{n_{1.}}{n}, \frac{n_{2.}}{n}, \dots, \frac{n_{r.}}{n} \right)$$

e

$$\bar{p}^c = \left(\frac{n_{.1}}{n}, \frac{n_{.2}}{n}, \dots, \frac{n_{.c}}{n} \right)$$

representam os perfis médios de Linhas e Colunas.



Exemplo Viagens (Greenacre, 2008)

Programa ACGreenacre08.R - EV.Dados

Tabela: Número de viagens em relação a atividade e localidade.

País	Descanso	Trabalho	Trabalho	Total
		Parcial	Integral	
Noruega	6	1	11	18
Canadá	1	3	11	15
Grécia	4	25	0	29
França/Alemanha	2	2	20	24
Total	13	31	42	86



Exemplo Viagens (Greenacre, 2008) - Perfis e Espaço de Perfis

Programa ACGreenacre08.R - EV.Perfis e Espaço Perfis

Tabela: Perfis de Linhas.

País	Descanso	Trabalho Parcial	Trabalho Integral	Total
Noruega	0.33	0.06	0.61	1.00
Canadá	0.07	0.20	0.73	1.00
Grécia	0.14	0.86	0.00	1.00
França/Alemanha	0.08	0.08	0.83	1.00
Média	0.15	0.36	0.49	1.00

Tabela: Perfis de Colunas.

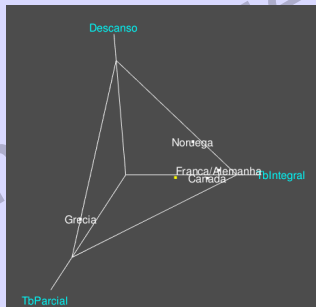
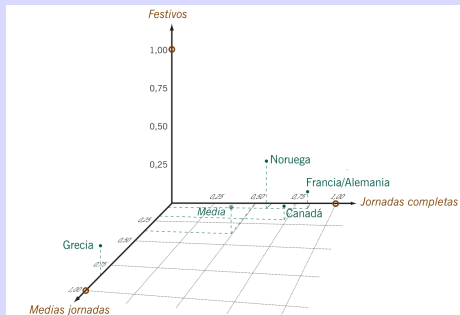
País	Descanso	TbParcial	TbIntegral	Média
Noruega	0.46	0.03	0.26	0.21
Canadá	0.08	0.10	0.26	0.17
Grécia	0.31	0.81	0.00	0.34
França/Alemanha	0.15	0.06	0.48	0.28
Total	1.00	1.00	1.00	1.00



Exemplo Viagens (Greenacre, 2008) Perfis e Espaço de Perfis

Programa ACGreenacre08.R - EV.Gráfico ($r = 3$)

Perfis de Linha (Países) no Triângulo de Colunas (Atividade)



Greenacre, 2008

Código R

- Mas como representar perfis de uma tabela com mais de três elementos nas linhas e colunas? Esta é a situação da Tabela de Izenman (Slide 3), que tem 4 linhas (cores dos olhos) e 5 colunas (cores dos cabelos).



Exemplo Leitura (Greenacre, 2008)

Programa ACGreenacre08.R - EL.Dados e EL.Massas e Centroides

Massas e Centróides

Tabela: Tipo de Leitor e Nível Educacional (Greenacre, 2008)

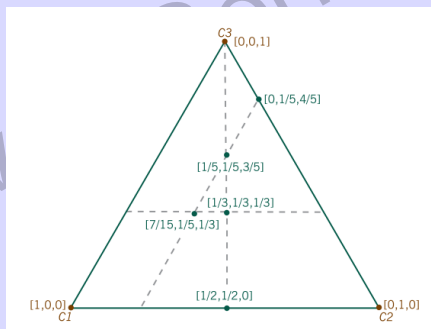
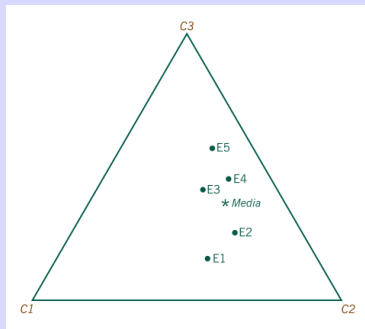
NIVEL DE EDUCACIÓN	TIPO DE LECTOR			Total	Masas de las filas
	Rápidos C1	Minuciosos C2	Muy minuciosos C3		
Educación primaria incompleta E1	5 (0,357)	7 (0,500)	2 (0,143)	14	0,045
Educación primaria E2	18 (0,214)	46 (0,548)	20 (0,238)	84	0,269
Educación secundaria incompleta E3	19 (0,218)	29 (0,333)	39 (0,448)	87	0,279
Educación secundaria E4	12 (0,119)	40 (0,396)	49 (0,485)	101	0,324
Educación universitaria incompleta E5	3 (0,115)	7 (0,269)	16 (0,615)	26	0,083
Total	57	129	126	312	
Perfil fila medio	(0,183)	(0,413)	(0,404)		



Exemplo Leitura (Greenacre, 2008)

Programa ACGreenacre08.R - EL.Gráfico ($r = 3$)

Massas e Centróides - Representação Gráfica



Exemplo Leitura (Greenacre, 2008)

Programa ACGreenacre08.R - EL.Freq Esperadas

Tabela: Frequências Observadas e Esperadas* (entre parênteses)

NIVEL DE EDUCACIÓN	TIPO DE LECTOR			Total	Masas de las filas
	Rápidos C1	Minuciosos C2	Muy minuciosos C3		
Educación primaria incompleta E1	5 (2,56)	7 (5,78)	2 (5,66)	14	0,045
Educación primaria E2	18 (15,37)	46 (34,69)	20 (33,94)	84	0,269
Educación secundaria incompleta E3	19 (15,92)	29 (35,93)	39 (35,15)	87	0,279
Educación secundaria E4	12 (18,48)	40 (41,71)	49 (40,80)	101	0,324
Educación universitaria incompleta E5	3 (4,76)	7 (10,74)	16 (10,50)	26	0,083
Total	57	129	126	312	
Perfil fila medio	(0,183)	(0,413)	(0,404)		

* Sob suposição de homogeneidade/independência. Exemplo: $2,56 = 14 \times 57/312$.



Distância Qui-Quadrado e Inércia

Programa ACGreenacre08.R - EL.Qui-Quadrado

- A estatística χ^2 permite testar as hipóteses

$$\begin{cases} H_0 : \text{Dados Homogêneos} \\ H_1 : \text{Dados Não Homogêneos} \end{cases}$$

- Estatística

$$\chi_0^2 = \sum \frac{(\text{observado} - \text{esperado})^2}{\text{esperado}} = \sum_{j=1}^c \sum_{i=1}^r \frac{(O_{ij} - E_{ij})^2}{E_{ij}}$$

tem distribuição $\chi_{(r-1)(c-1)}^2$.

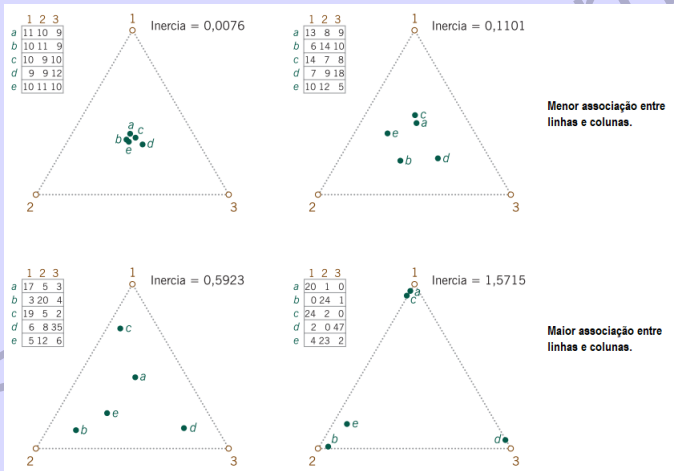
- Para a tabela em questão (leitor) teremos $\chi_0^2 = 25.98$ e $P(\chi_8^2 > 25.98) = 0.00106$. Logo, rejeitamos a hipótese de homogeneidade nos dados, i.e., existe uma associação entre o nível de educação e tipo de leitor observados.



Distância Qui-Quadrado e Inércia

Programa ACGreenacre08.R - EL.Qui-Quadrado

- A Inércia (in) total é igual a χ^2/n e mede a variação dos perfis individuais \mathbf{p}_i (ou \mathbf{p}_j) em torno do centróide $\bar{\mathbf{p}}$.
- Para os dados em questão temos $in = 0.08326$. Significado?



Análise de Correspondências

Abordagem Matricial



Programa ACMatricial.R - NBasica

- Considere uma tabela $\mathbf{M}_{(2 \times 3)}$ de contingência artificial:

	B ₁	B ₂	B ₃	Total
A ₁	1	3	2	$n_{1.} = 6$
A ₂	2	2	3	$n_{2.} = 7$
Total	$n_{.1} = 3$	$n_{.2} = 5$	$n_{.3} = 5$	$n = 13$

- Defina,

$$\mathbf{X}_{ij} = \begin{cases} 1 & \text{se indivíduo } j \text{ pertence a } A_i, j = 1, \dots, c \\ 0 & \text{caso contrário.} \end{cases}$$

$$\mathbf{Y}_{ij} = \begin{cases} 1 & \text{se indivíduo } i \text{ pertence a } B_j, i = 1, \dots, r \\ 0 & \text{caso contrário.} \end{cases}$$

- Então,

$$\mathbf{X} = \left[\begin{array}{c|ccc|cc|cc|cc|ccc} 1 & 1 & 1 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \end{array} \right]$$

$$\mathbf{Y} = \left[\begin{array}{c|ccc|cc|cc|cc|ccc} 1 & 0 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 1 & 0 & 0 & 0 & 0 & 1 & 1 & 1 \end{array} \right]$$



Seguindo esta notação,

- $\mathbf{XY}^T = \mathbf{M}$;
- $\mathbf{XX}^T = \text{diag}(n_{1.}, \dots, n_{r.})$;
- $\mathbf{YY}^T = \text{diag}(n_{.1}, \dots, n_{.c})$;
- $\mathbf{D}_r = \mathbf{XX}^T / n = \text{diag}(n_{1.}/n, \dots, n_{r.}/n)$;
- $\mathbf{D}_c = \mathbf{YY}^T / n = \text{diag}(n_{.1}/n, \dots, n_{.c}/n)$ e
- Matriz de Burt,

$$\mathbf{B} = \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix} \begin{bmatrix} \mathbf{X} \\ \mathbf{Y} \end{bmatrix}^T = \begin{bmatrix} n \mathbf{D}_r & \mathbf{M} \\ \mathbf{M}^T & n \mathbf{D}_c \end{bmatrix}$$

positiva semi-definida e simétrica.



Perfis, Massa e Centróides

- A matriz de correspondências \mathbf{P} é obtida por $\mathbf{XY}^T/n = \mathbf{M}/n$;
- Os totais das linhas e colunas de \mathbf{P} são as diagonais de \mathbf{D}_r e \mathbf{D}_c .
- Os Perfis das linhas de \mathbf{M} são obtidos por $\mathbf{P}_r = \mathbf{D}_r^{-1}\mathbf{P} = (\mathbf{XX}^T)^{-1}\mathbf{XY}^T$, i.e.,
corresponde aos coeficientes da regressão de \mathbf{Y} em \mathbf{X} ;
- Os Perfis das colunas de \mathbf{M} são obtidos por $\mathbf{P}_c = \mathbf{D}_c^{-1}\mathbf{P}^T = (\mathbf{YY}^T)^{-1}\mathbf{YX}^T$, i.e.,
corresponde aos coeficientes da regressão de \mathbf{X} em \mathbf{Y} ;
- Seja \mathbf{r} o vetor formado pelos elementos da diagonal de \mathbf{D}_r . Então $\mathbf{r} = \mathbf{P}\mathbf{1}_c$ representa as proporções dos totais de linhas em relação ao total. Também conhecido como centróide (média de perfis) de colunas.
- Seja \mathbf{c} o vetor formado pelos elementos da diagonal de \mathbf{D}_c . Então $\mathbf{c} = \mathbf{P}^T\mathbf{1}_r$ representa as proporções dos totais de colunas em relação ao total. Também conhecido como centróide de perfis de linhas.



Distâncias Qui-Quadrado

- Seja \mathbf{a}_i o i -ésimo perfil de linha, $i = 1, \dots, r$, i.e., \mathbf{a}_i corresponde a i -ésima linha de \mathbf{P}_r .
- Seja \mathbf{b}_j o j -ésimo perfil de coluna, $j = 1, \dots, c$. \mathbf{b}_j corresponde a j -ésima linha de \mathbf{P}_c .
- A distância χ^2 entre \mathbf{a}_i e $\mathbf{a}_{i'}$ é definida pela forma Quadrática:

$$\mathbf{d}^2 = (\mathbf{a}_i - \mathbf{a}_{i'})^T \mathbf{D}_c^{-1} (\mathbf{a}_i - \mathbf{a}_{i'}) = \sum_{j=1}^c \frac{n}{n_{.j}} \left(\frac{n_{ij}}{n_{i.}} - \frac{n_{i'j}}{n_{i'.}} \right)^2.$$

$n/n_{.j}$ indica que as colunas com menos elementos contribuem mais para as distâncias entre perfis de linhas.

- \mathbf{c} é o centróide de colunas. A matrix $(r \times c)$ de perfis centralizados, $\mathbf{P}_r - \mathbf{1}_r \mathbf{c}^T$ tem i -ésima linha igual a $(\mathbf{a}_i - \mathbf{c})^T$, com j -ésimo elemento

$$\frac{1}{n_{i.}} \left(n_{ij} - \frac{n_{i.} n_{.j}}{n} \right) \quad i = 1, \dots, r; \quad j = 1, \dots, c.$$



Distâncias Qui-Quadrado - continuação

- Assim, a distância χ^2 entre \mathbf{a}_i e \mathbf{c} é

$$d^2(\mathbf{a}_i, \mathbf{c}) = (\mathbf{a}_i - \mathbf{c})^T \mathbf{D}_c^{-1} (\mathbf{a}_i - \mathbf{c}) = \frac{1}{n_i} \sum_{j=1}^c \frac{n}{n_j \cdot n_j} \left(n_{ij} - \frac{n_i \cdot n_j}{n} \right)^2 \quad (1)$$

- Somando a Equação 1 para todo perfil de linha fornece

$$n \sum_{i=1}^r p_i \cdot d^2(\mathbf{a}_i, \mathbf{c}) = \sum_{i=1}^r \sum_{j=1}^c \left(n_{ij} - \frac{n_i \cdot n_j}{n} \right)^2 / \left(\frac{n_i \cdot n_j}{n} \right),$$

que é a estatística χ^2 , uma vez que $O_{ij} = n_{ij}$ e $E_{ij} = \frac{n_i \cdot n_j}{n}$. Já vimos que para amostragem aleatória $\chi^2 \sim \chi^2_{(r-1)(c-1)}$.

- Então $\chi^2_0 = 1240.04$ e $P(\chi^2_{12} > 1240.04) \approx 0$, e concluímos que existe associação significativa entre linhas e colunas.
- O mesmo raciocínio pode ser feito utilizando \mathbf{b}_j .



Inércia Total e Decomposição - Resultados

- Sejam $\mathbf{x} = (x_{ij})$, em que $x_{ij} = (X_{ij} - \bar{X}_i)$ e $\mathbf{y} = (y_{ij})$, em que $y_{ij} = (Y_{ij} - \bar{Y}_j)$. As matrizes de covariância são

$$n^{-1}\mathbf{xx}^T = \mathbf{D}_r - \mathbf{r}\mathbf{r}^T \quad \text{e} \quad n^{-1}\mathbf{yy}^T = \mathbf{D}_c - \mathbf{c}\mathbf{c}^T.$$

Estas matrizes são singulares e portanto não tem inversa.

- A matriz de resíduos pode ser obtida por

$$\tilde{\mathbf{P}} = \mathbf{P} - \mathbf{r}\mathbf{c}^T$$

por que seus elementos são iguais a $\tilde{n}_{ij} = O_{ij} - E_{ij}$.



Inércia Total e Decomposição - continuação

- A distância χ^2 pode ser obtida também por

$$\mathbf{R}_0 = \mathbf{D}_c^{-1/2} \tilde{\mathbf{P}}^T \mathbf{D}_r^{-1} \tilde{\mathbf{P}} \mathbf{D}_c^{-1/2},$$

uma vez que pode-se mostrar que $\text{tr}(\mathbf{R}_0) = \frac{\chi^2}{n} = \text{Inércia Total (in)}$ da Tabela de Contingência.

- A Inércia Total para os dados em questão é igual a 0.23, indicando uma relação razoável em cor dos olhos e cor dos cabelos.
- Assim como em PCA, os autovalores de \mathbf{R}_0 indicam a contribuição das $t = \max(i, j)$ inércias principais na inércia total. No caso em questão, as duas primeiras inércias principais contribuem para cerca de 99.63% da inérica total. Isto indica que a matriz de dados pode ser bem representada no plano.



Coordenadas Principais para Perfis de Linhas e Colunas

- \mathbf{R}_0 pode ser representada por $\mathbf{L} = \mathbf{D}_r^{-1/2} \tilde{\mathbf{P}} \mathbf{D}_c^{-1/2}$ que tem como ℓ_{ij} o resíduo de Pearson, i.e.,

$$\ell_{ij} = (n_{i.} n_{.j})^{-1/2} \left(n_{ij} - \frac{n_{i.} n_{.j}}{n} \right) \quad i = 1, \dots, r; \quad j = 1, \dots, c.$$

- Assim, a inércia também pode ser obtida por $\sum_i \sum_j \ell_{ij}^2$.
- Resultados:
 - $\mathbf{R}_0 = \mathbf{L}^T \mathbf{L}$.
 - Seja $\mathbf{R}_1 = \mathbf{L} \mathbf{L}^T$.
 - A matriz \mathbf{L} pode ser escrita pela decomposição em valores singulares $\mathbf{U} \mathbf{D}_\lambda \mathbf{V}^T$. Os autovetores de \mathbf{R}_1 são iguais a \mathbf{U} e os autovetores de \mathbf{R}_0 são iguais as colunas de \mathbf{V} .
 - Note que

$$\tilde{\mathbf{P}} = \mathbf{D}_r^{1/2} \mathbf{L} \mathbf{D}_c^{1/2} = (\mathbf{D}_r^{1/2} \mathbf{U}) \mathbf{D}_\lambda (\mathbf{V}^T \mathbf{D}_c^{1/2}) = \mathbf{A} \mathbf{D}_\lambda \mathbf{B}^T$$



Coordenadas Principais para Perfis de Linhas e Colunas - continuação

- Ainda, $\mathbf{A}^T \mathbf{D}_r^{-1} \mathbf{A} = \mathbf{I}_c$ e $\mathbf{B}^T \mathbf{D}_c^{-1} \mathbf{B} = \mathbf{I}_c$.
- A expressão $\tilde{\mathbf{P}} = \mathbf{D}_r^{1/2} \mathbf{L} \mathbf{D}_c^{1/2} = (\mathbf{D}_r^{1/2} \mathbf{U}) \mathbf{D}_\lambda (\mathbf{V}^T \mathbf{D}_c^{1/2}) = \mathbf{A} \mathbf{D}_\lambda \mathbf{B}^T$ é a decomposição em valores singulares de $\tilde{\mathbf{P}}$ nas métricas \mathbf{D}_r^{-1} e \mathbf{D}_c^{-1} .
- As colunas de \mathbf{A} são chamadas de eixos principais de linhas.
- As colunas de \mathbf{B} são chamadas de eixos principais de colunas.
- As coordenadas principais dos perfis de linhas (\mathbf{G}_P^T) é a distância χ^2 entre a matriz de perfis centralizados de linhas $(\mathbf{P}_r - \mathbf{1}_r \mathbf{c}^T)$ e a matriz \mathbf{B} , na métrica¹ \mathbf{D}_c^{-1} , i.e.,

$$\mathbf{G}_P^T = (\mathbf{P}_r - \mathbf{1}_r \mathbf{c}^T) \mathbf{D}_c^{-1} \mathbf{B} = \mathbf{D}_r^{-1} \mathbf{A} \mathbf{D}_\lambda.$$
- As coordenadas principais dos perfis de colunas (\mathbf{H}_P^T) é a distância χ^2 entre a matriz de perfis centralizados de colunas $(\mathbf{P}_c - \mathbf{1}_c \mathbf{r}^T)$ e a matriz \mathbf{A} , na métrica² \mathbf{D}_r^{-1} , i.e.,

$$\mathbf{H}_P^T = (\mathbf{P}_c - \mathbf{1}_c \mathbf{r}^T) \mathbf{D}_r^{-1} \mathbf{A} = \mathbf{D}_c^{-1} \mathbf{B} \mathbf{D}_\lambda.$$

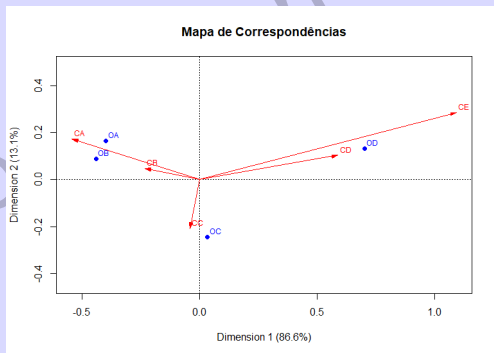
¹Lembre que \mathbf{D}_c é uma matriz diagonal com os pesos dos totais de colunas da tabela em análise.

²Lembre que \mathbf{D}_r é uma matriz diagonal com os pesos dos totais de linhas da tabela em análise.

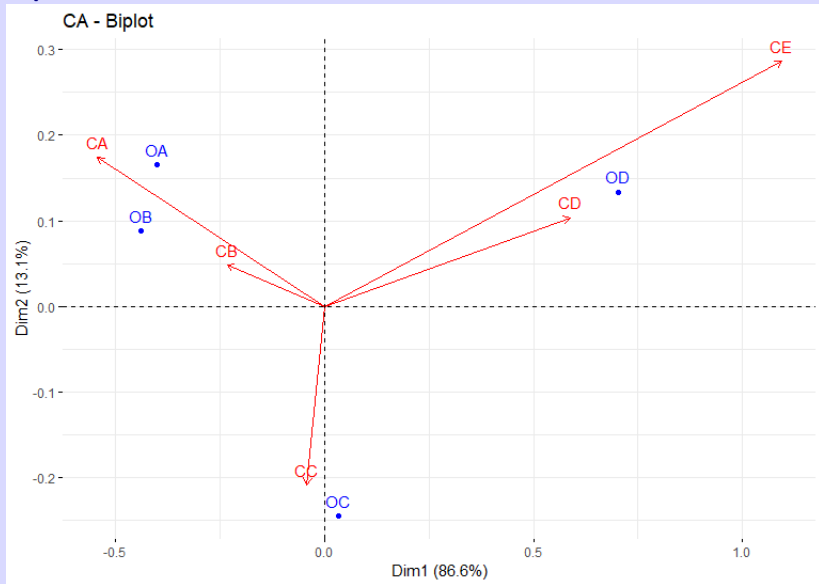


Mapas de Correspondências

- **Mapa Simétrico:** as coordenadas de linhas e colunas podem ser expressas como coordenadas principais.
- **Mapa Assimétrico:** As coordenadas de linhas (ou colunas) são expressas como coordenadas principais e a outra coordenada como coordenadas padrão.
- A coordenada padrão de linhas é $\mathbf{G}_s = \mathbf{U}^T \mathbf{D}_r^{-1/2}$.
- A coordenada padrão de colunas é $\mathbf{H}_s = \mathbf{V}^T \mathbf{D}_c^{-1/2}$.



Biplot com Pacote FactoMiner



Biplot com Pacote FactoMiner

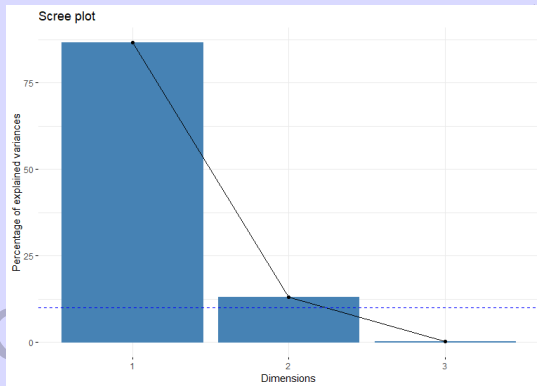
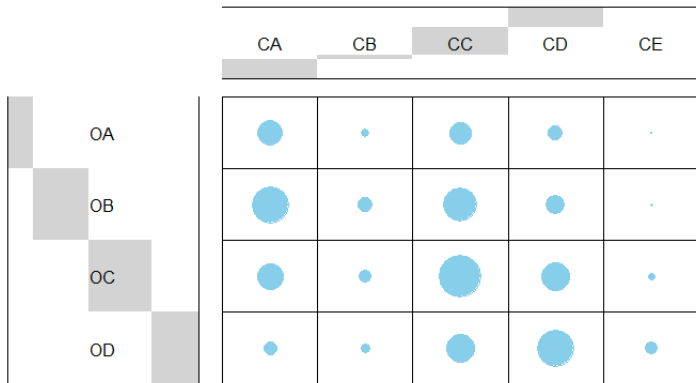


Grafico com gplots

Rel. Cor dos Olhos e Cabelos



CA em Tabelas Multidimensionais

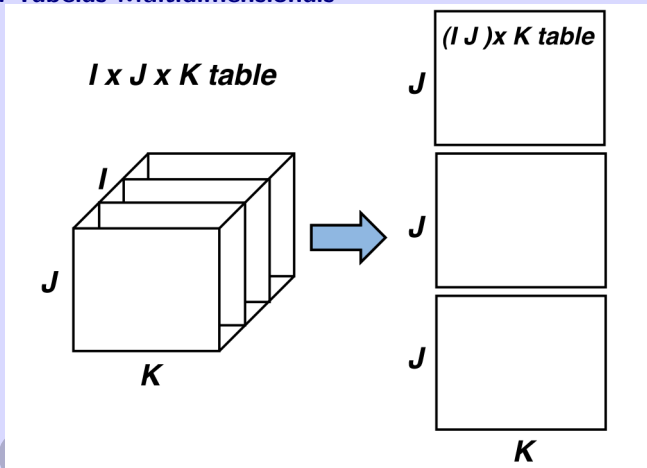


Figura: Unindo três carecterísticas numa tabela bidimensional.

Fonte: Friendly e Meyer (2016).

