

Análise de Correspondência

Prof. Caio Azevedo

- As metodologias de ACP (análise de componentes principais) e de AF (análise fatorial) nos permitem, entre outros aspectos:
 - Compreender melhor o comportamento das unidades amostrais em relação às variáveis originais.
 - Compreender melhor a estrutura de dependência entre as variáveis originais.
- Discutimos, brevemente, que a análise de dados categorizados também faz parte da análise de dados multivariados.
- As metodologias de ACP e de AF não são apropriadas para analisar dados categorizados.

Exemplo 6: dados sobre aquisição de aparelhos de som

- Foram entrevistados 1320 consumidores de aparelhos de som.
- Variáveis medidas: marca adquirida (cinco categorias) e principal motivo da compra (seis categorias).
- Ou seja, o modelo probabilístico gerador da tabela de contingência em questão é uma multinomial de tamanho 1320 com 30 categorias.
- Os dados obtidos encontram-se na tabela a seguir:

Marca	Atributo						Total
	Qual. de som	Tec.	Pot. do som	Rec. técnicos	Preço	Conf. na marca	
Sony	135	140	95	55	40	60	525
Aiwa	50	115	40	60	5	15	285
Gradiente	90	55	20	35	40	10	250
Philips	60	25	35	10	5	30	165
Sharp	30	20	5	10	10	20	95
Total	365	355	195	170	100	135	1320

OBS: Qual. de som - qualidade do som; tec. - tecnologia avançada; pot. do som - potência do som; rec. técnicos - recursos técnicos ; conf. na marca - confiança na marca.

- Formal geral de tabelas de contingência (2×2).

Variável 1	Variável 2					Total
	Cat. 1	Cat. 2	Cat. 3	...	Cat. J	
Categoria 1	X_{11}	X_{12}	X_{13}	...	X_{1J}	$X_{1.}$
Categoria 2	X_{21}	X_{22}	X_{23}	...	X_{2J}	$X_{2.}$
Categoria 3	X_{31}	X_{32}	X_{33}	...	X_{3J}	$X_{3.}$
	\vdots	\vdots	\vdots	\ddots	\vdots	\vdots
Categoria I	X_{I1}	X_{I2}	X_{I3}	...	X_{IJ}	$X_{I.}$
Total	$X_{.1}$	$X_{.2}$	$X_{.3}$...	$X_{.J}$	$X_{..}$

- Seja $Y_{ijk} = 1$ se o indivíduo k , $k=1,2,\dots,n$, foi classificado na categoria i da variável 1 e na categoria j da variável 2, tal que $\mathbf{Y}_k = (Y_{11k}, \dots, Y_{IJk})' \stackrel{i.i.d.}{\sim} \text{multinomial}(n, \mathbf{p})$, $\mathbf{p} = (p_{11}, \dots, p_{IJ})'$. Assim $X_{ij} = \sum_{k=1}^n Y_{ijk}$.
- Seja $p_{ij} = P(X_{ij} = 1)$.
- Sob algumas suposições e para $n = X_{..}$ fixado, temos que o número de unidades amostrais observadas nas $I \times J$ categorias seguem uma (conjuntamente) uma distribuição multinomial n , p_{ij} .

- Hipótese de interesse: H_0 : as variáveis são estatisticamente independentes vs H_1 : as variáveis não são estatisticamente independentes.
- $H_0 : p_{ij} = p_{i.}p_{.j}, \forall i, j$ vs $H_1 : p_{ij} \neq p_{i.}p_{.j}$, para pelo menos um par (i, j) .
- $p_{i.}$ e $p_{.j}$ são as probabilidades marginais de cada unidade amostral pertencer, respectivamente à categoria i da variável 1 e à categoria j da variável 2.

- $p_{i.} = \sum_{j=1}^J p_{ij}$ e $p_{.j} = \sum_{i=1}^I p_{ij}$.
- Sob independência, temos que $Q = \sum_{i=1}^I \sum_{j=1}^J \frac{(X_{ij} - E_{ij})^2}{E_{ij}}$, em que $E_{ij} = X_{i.} X_{.j} / X_{..}$, tenderá a apresentar um valor baixo.
- E_{ij} é a frequência (valor esperado) da casela (i, j) sob independência.
- Sob a validade de suposição de independência, e para tamanhos amostrais suficientemente grandes, temos que $Q \approx \chi_{(I-1)(J-1)}^2$.

- Verificar se existe associação entre marca e atributo, em termos de aquisição, por parte dos consumidores.
- Estatística de qui-quadrado: $Q = 179,62$, $p\text{-valor} < 0,0001$.
- Portanto existe dependência entre as duas variáveis.
- Investigar, minuciosamente, a relação de dependência entre as variáveis que definem a tabela de contingência.
- Por exemplo: que marca de aparelho de som é mais adquirida em função do preço?

- Análise de correspondência: visa medir o grau de associação de variáveis categorizadas dispostas em tabelas de contingência.
- A disposição dos resultados é feita de modo gráfico.
- Uma forma de medir associação é através da estatística de qui-quadrado.
- Alternativa: modelos lineares, log-lineares e não lineares para dados categorizados. Veja mais em:

http:

[//www.ime.unicamp.br/~cnaber/Material_ADD_1S_2017.htm](http://www.ime.unicamp.br/~cnaber/Material_ADD_1S_2017.htm)

- Defina

$$\mathbf{X} = \begin{bmatrix} X_{11} & X_{12} & \dots & X_{1J} \\ X_{21} & X_{22} & \dots & X_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ X_{I1} & X_{I2} & \dots & X_{IJ} \end{bmatrix}; \mathbf{E} = \begin{bmatrix} E_{11} & E_{12} & \dots & E_{1J} \\ E_{21} & E_{22} & \dots & E_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ E_{I1} & E_{I2} & \dots & E_{IJ} \end{bmatrix}$$

- $\mathbf{X}^* = \text{vec}(\mathbf{X}') = (X_{11}, X_{12}, \dots, X_{IJ})'$ (frequências observadas) e
 $\mathbf{E}^* = \text{vec}(\mathbf{E}') = (E_{11}, E_{12}, \dots, E_{IJ})'$ (frequências esperadas sob independência).
- $X_{i.} = \sum_{j=1}^J X_{ij}$, $X_{.j} = \sum_{i=1}^I X_{ij}$ e $X_{..} = \sum_{i=1}^I \sum_{j=1}^J X_{ij}$.

- Note que $Q = (\mathbf{X}^* - \mathbf{E}^*)' \mathbf{D}_{\mathbf{E}^*}^{-1} (\mathbf{X}^* - \mathbf{E}^*)$ (estatística de Pearson).

$$\mathbf{D}_{\mathbf{E}^*} = \begin{bmatrix} E_{11} & 0 & \dots & 0 \\ 0 & E_{12} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & E_{IJ} \end{bmatrix}$$

- Defina $\mathbf{P} = \mathbf{X}/n = \mathbf{X}/X_{..}$ e $\mathbf{P}_E = \mathbf{E}/X_{..}$, e:

$$\mathbf{P} = \begin{bmatrix} P_{11} & P_{12} & \dots & P_{1J} \\ P_{21} & P_{22} & \dots & P_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ P_{I1} & P_{I2} & \dots & P_{IJ} \end{bmatrix}; \mathbf{P}_E = \begin{bmatrix} P_{E11} & P_{E12} & \dots & P_{E1J} \\ P_{E21} & P_{E22} & \dots & P_{E2J} \\ \vdots & \vdots & \ddots & \vdots \\ P_{EI1} & P_{EI2} & \dots & P_{EIJ} \end{bmatrix}$$

- Pode-se demonstrar que $Q = n(\mathbf{P}^* - \mathbf{P}_E^*)' \mathbf{D}_{\mathbf{P}_E}^{-1} (\mathbf{P}^* - \mathbf{P}_E^*)$ (exercício), em que $\mathbf{P}^* = \text{vec}(\mathbf{P}')$ e $\mathbf{P}_E^* = \text{vec}(\mathbf{P}'_E)$.

$$\mathbf{D}_{\mathbf{P}_E} = \begin{bmatrix} P_{E_{11}} & 0 & \dots & 0 \\ 0 & P_{E_{12}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & P_{E_{JJ}} \end{bmatrix}$$

- A análise de correspondência explora a forma acima.
- Mais especificamente, explora o que chamamos de: perfil das linhas e perfil das colunas.

- Perfil das linhas: $p_{j/i} = \frac{p_{ij}}{p_{i.}}$ (probabilidade de pertencer à categoria j da variável 2 dado que pertence à categoria i da variável 1).
- Perfil das colunas: $p_{i/j} = \frac{p_{ij}}{p_{.j}}$ (probabilidade de pertencer à categoria i da variável 1 dado que pertence à categoria j da variável 2).
- Definamos:
 - $P_r = P\mathbf{1}_J$.
 - $P_c = P'\mathbf{1}_I$.

$$\begin{aligned}
 \blacksquare \mathbf{P}_{r(I \times 1)} &= \begin{bmatrix} P_{11} & P_{12} & \dots & P_{1J} \\ P_{21} & P_{22} & \dots & P_{2J} \\ \vdots & \vdots & \ddots & \vdots \\ P_{I1} & P_{I2} & \dots & P_{IJ} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} \sum_{j=1}^J P_{1j} \\ \sum_{j=1}^J P_{2j} \\ \vdots \\ \sum_{j=1}^J P_{Ij} \end{bmatrix} = \\
 \begin{bmatrix} P_{1.} \\ P_{2.} \\ \vdots \\ P_{I.} \end{bmatrix} &= \begin{bmatrix} \frac{X_{1.}}{X_{..}} \\ \frac{X_{2.}}{X_{..}} \\ \vdots \\ \frac{X_{I.}}{X_{..}} \end{bmatrix}.
 \end{aligned}$$

$$\begin{aligned}
 \blacksquare \mathbf{P}_{c(J \times 1)} &= \begin{bmatrix} P_{11} & P_{21} & \dots & P_{I1} \\ P_{12} & P_{22} & \dots & P_{I2} \\ \vdots & \vdots & \ddots & \vdots \\ P_{1J} & P_{2J} & \dots & P_{IJ} \end{bmatrix} \begin{bmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{bmatrix} = \begin{bmatrix} \sum_{i=1}^I P_{i1} \\ \sum_{i=1}^I P_{i2} \\ \vdots \\ \sum_{i=1}^I P_{iJ} \end{bmatrix} = \\
 \begin{bmatrix} P_{.1} \\ P_{.2} \\ \vdots \\ P_{.J} \end{bmatrix} &= \begin{bmatrix} \frac{X_{.1}}{X_{..}} \\ \frac{X_{.2}}{X_{..}} \\ \vdots \\ \frac{X_{.J}}{X_{..}} \end{bmatrix}.
 \end{aligned}$$

- Perfis das linhas (forma matricial):

$$R = D_r^{-1}P = \begin{bmatrix} \frac{P_{11}}{P_{1.}} & \frac{P_{12}}{P_{1.}} & \cdots & \frac{P_{1J}}{P_{1.}} \\ \frac{P_{21}}{P_{2.}} & \frac{P_{22}}{P_{2.}} & \cdots & \frac{P_{2J}}{P_{2.}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{P_{I1}}{P_{I.}} & \frac{P_{I2}}{P_{I.}} & \cdots & \frac{P_{IJ}}{P_{I.}} \end{bmatrix} = \begin{bmatrix} \frac{X_{11}}{X_{1.}} & \frac{X_{12}}{X_{1.}} & \cdots & \frac{X_{1J}}{X_{1.}} \\ \frac{X_{21}}{X_{2.}} & \frac{X_{22}}{X_{2.}} & \cdots & \frac{X_{2J}}{X_{2.}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{I1}}{X_{I.}} & \frac{X_{I2}}{X_{I.}} & \cdots & \frac{X_{IJ}}{X_{I.}} \end{bmatrix}$$

$$D_r = \begin{bmatrix} P_{1.} & 0 & \cdots & 0 \\ 0 & P_{2.} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & P_{I.} \end{bmatrix}$$

- Perfis das colunas:

$$C = D_c^{-1}P' = \begin{bmatrix} \frac{P_{11}}{P_{.1}} & \frac{P_{21}}{P_{.1}} & \cdots & \frac{P_{J1}}{P_{.1}} \\ \frac{P_{12}}{P_{.2}} & \frac{P_{22}}{P_{.2}} & \cdots & \frac{P_{J2}}{P_{.2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{P_{1J}}{P_{.J}} & \frac{P_{2J}}{P_{.J}} & \cdots & \frac{P_{JJ}}{P_{.J}} \end{bmatrix} = \begin{bmatrix} \frac{X_{11}}{X_{.1}} & \frac{X_{21}}{X_{.1}} & \cdots & \frac{X_{J1}}{X_{.1}} \\ \frac{X_{12}}{X_{.2}} & \frac{X_{22}}{X_{.2}} & \cdots & \frac{X_{J2}}{X_{.2}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{1J}}{X_{.J}} & \frac{X_{2J}}{X_{.J}} & \cdots & \frac{X_{JJ}}{X_{.J}} \end{bmatrix}$$

$$D_c = \begin{bmatrix} P_{.1} & 0 & \cdots & 0 \\ 0 & P_{.2} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & P_{.J} \end{bmatrix}$$

- Proporções estimadas (pelo total de observações) ($\times 100$).

Marca	Atributo						Total
	Qual. de som	Tec.	Pot. do som	Rec. técnicos	Preços	Conf. na marca	
Sony	10,23	10,61	7,20	4,17	3,03	4,55	39,77
Aiwa	3,79	8,71	3,03	4,55	0,38	1,14	21,59
Gradiente	6,82	4,17	1,52	2,65	3,03	0,76	18,94
Philips	4,55	1,89	2,65	0,76	0,38	2,27	12,50
Sharp	2,27	1,52	0,38	0,76	0,76	1,52	7,20
Total	27,65	26,89	14,77	12,88	7,58	10,23	100,00

- Ranqueamento dos atributos e das marcas.

Marca	Atributo						Rank
	Qual. de som	Tec.	Pot. do som	Rec. técnicos	Preços	Conf. na marca	
Sony	10,23	10,61	7,20	4,17	3,03	4,55	1
Aiwa	3,79	8,71	3,03	4,55	0,38	1,14	2
Gradiente	6,82	4,17	1,52	2,65	3,03	0,76	3
Philips	4,55	1,89	2,65	0,76	0,38	2,27	4
Sharp	2,27	1,52	0,38	0,76	0,76	1,52	5
Rank	1	2	3	4	6	5	1

■ Perfis das linhas ($\times 100$)

Marca	Atributo						Total
	Qual. de som	Tec.	Pot. do som	Rec. técnicos	Preços	Conf. na marca	
Sony	25,71	26,67	18,10	10,48	7,62	11,43	100
Aiwa	17,54	40,35	14,04	21,05	1,75	5,26	100
Gradiente	36,00	22,00	8,00	14,00	16,00	4,00	100
Philips	36,36	15,15	21,21	6,06	3,03	18,18	100
Sharp	31,58	21,05	5,26	10,53	10,53	21,05	100
Total	27,65	26,89	14,77	12,88	7,58	10,23	100

■ Perfis das colunas ($\times 100$)

Marca	Atributo						
	Qual. de som	Tec.	Pot. do som	Rec. técnicos	Preços	Conf. na marca	Total
Sony	36,99	39,44	48,72	32,35	40,00	44,44	39,77
Aiwa	13,70	32,39	20,51	35,29	5,00	11,11	21,59
Gradiente	24,66	15,49	10,26	20,59	40,00	7,41	18,94
Philips	16,44	7,04	17,95	5,88	5,00	22,22	12,50
Sharp	8,22	5,63	2,56	5,88	10,00	14,81	7,20
Total	100	100	100	100	100	100	100

- A observação das tabelas de perfis de linhas e colunas já pode fornecer alguma idéia das associações existentes.
- Na tabela dos perfis das linhas, pode-se perceber que, no total, 27,65% escolhem uma determinada marca em função da qualidade do som e 7,58% escolhem por causa do preço.
- Com relação aos consumidores que compraram o produto da marca Gradiente, 36,00% o fizeram devido à qualidade do som e 16,00% devido ao preço, que são maiores do que as respectivas proporções na população geral (27,65% e 7,58%).

- A marca Sony apresenta uma porcentagem maior do que a observada na população geral, destacadamente, apenas no atributo potência de som, indicando que esta marca pode estar mais relacionada com tal atributo.
- Na tabela dos perfis das colunas, nota-se que os consumidores preocupados com o atributo tecnologia avançada, 32,39% escolhem a marca Aiwa, que é maior do que 21,59% (porcentagem observada na população geral).

- Objetivo: obter uma forma simplificada e interpretável para as seguintes matrizes:

$$\mathbf{W}_{(I \times J)} = \mathbf{D}_r^{1/2} (\mathbf{R} - \mathbf{1P}'_c) \mathbf{D}_c^{-1/2}$$

$$\mathbf{Z}_{(J \times I)} = \mathbf{D}_c^{1/2} (\mathbf{C} - \mathbf{1P}'_r) \mathbf{D}_r^{-1/2}$$

- Simplificada: duas dimensões (gráfico de dispersão)
- Interpretável: pontos próximos (categorias) tem um maior grau de dependência.

- Lembrando que as matrizes \mathbf{R} e \mathbf{C} são as matrizes com os perfis das linhas e colunas, respectivamente.
- A matriz $(\mathbf{R} - \mathbf{1P}'_c)$ representa os perfis das linhas centrados. Isso porque o vetor \mathbf{P}_c representa o ponto médio das linhas.
- Note que o j -ésimo elemento de cada perfil das linhas é dado por $X_{ij}/X_{i.}$.
- O j -ésimo elemento do vetor que representa o ponto médio será igual a soma de todos os i -ésimos elementos dos perfis de linhas multiplicados por suas respectivas frequências relativas (de cada linha), ou seja : $\sum_{i=1}^I \left(\frac{X_{ij}}{X_{i.}} \frac{X_{i.}}{X_{..}} \right) = \frac{X_{.j}}{X_{..}}$.
- Analogamente, $(\mathbf{C} - \mathbf{1P}'_r)$ representam os perfis de colunas centrados.

- Note que:

$$\begin{aligned}
 R - \mathbf{1}P'_c &= \begin{bmatrix} \frac{X_{11}}{X_{1.}} & \frac{X_{12}}{X_{1.}} & \cdots & \frac{X_{1J}}{X_{1.}} \\ \frac{X_{21}}{X_{2.}} & \frac{X_{22}}{X_{2.}} & \cdots & \frac{X_{2J}}{X_{2.}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{J1}}{X_{J.}} & \frac{X_{J2}}{X_{J.}} & \cdots & \frac{X_{JJ}}{X_{J.}} \end{bmatrix} - \begin{bmatrix} \frac{X_{.1}}{X_{..}} & \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{.J}}{X_{..}} \\ \frac{X_{.1}}{X_{..}} & \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{.J}}{X_{..}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{.1}}{X_{..}} & \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{.J}}{X_{..}} \end{bmatrix} \\
 &= \begin{bmatrix} \frac{X_{11}}{X_{1.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{12}}{X_{1.}} - \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{1J}}{X_{1.}} - \frac{X_{.J}}{X_{..}} \\ \frac{X_{21}}{X_{2.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{22}}{X_{2.}} - \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{2J}}{X_{2.}} - \frac{X_{.J}}{X_{..}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{J1}}{X_{J.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{J2}}{X_{J.}} - \frac{X_{.2}}{X_{..}} & \cdots & \frac{X_{JJ}}{X_{J.}} - \frac{X_{.J}}{X_{..}} \end{bmatrix}
 \end{aligned}$$

- Por outro lado, temos que:

$$D_r = \begin{bmatrix} P_{1.} & 0 & \dots & 0 \\ 0 & P_{2.} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & P_{I.} \end{bmatrix} = \begin{bmatrix} \frac{X_{1.}}{X_{..}} & 0 & \dots & 0 \\ 0 & \frac{X_{2.}}{X_{..}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{X_{I.}}{X_{..}} \end{bmatrix}$$

- Portanto, temos que:

$$D_r^{1/2} (R - \mathbf{1}P'_c) = \begin{bmatrix} \sqrt{\frac{X_{1.}}{X_{..}}} & 0 & \dots & 0 \\ 0 & \sqrt{\frac{X_{2.}}{X_{..}}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\frac{X_{I.}}{X_{..}}} \end{bmatrix} \\ \times \begin{bmatrix} \frac{X_{11}}{X_{1.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{12}}{X_{1.}} - \frac{X_{.2}}{X_{..}} & \dots & \frac{X_{1J}}{X_{1.}} - \frac{X_{.J}}{X_{..}} \\ \frac{X_{21}}{X_{2.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{22}}{X_{2.}} - \frac{X_{.2}}{X_{..}} & \dots & \frac{X_{2J}}{X_{2.}} - \frac{X_{.J}}{X_{..}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{X_{I1}}{X_{I.}} - \frac{X_{.1}}{X_{..}} & \frac{X_{I2}}{X_{I.}} - \frac{X_{.2}}{X_{..}} & \dots & \frac{X_{IJ}}{X_{I.}} - \frac{X_{.J}}{X_{..}} \end{bmatrix}$$

■ Logo:

$$D_r^{1/2} (R - \mathbf{1}P'_c) = \begin{bmatrix} \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{11}}{X_{1.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{12}}{X_{1.}} - \frac{X_{.2}}{X_{..}} \right) & \cdots & \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{1J}}{X_{1.}} - \frac{X_{.J}}{X_{..}} \right) \\ \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{21}}{X_{2.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{22}}{X_{2.}} - \frac{X_{.2}}{X_{..}} \right) & \cdots & \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{2J}}{X_{2.}} - \frac{X_{.J}}{X_{..}} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{I1}}{X_{I.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{I2}}{X_{I.}} - \frac{X_{.2}}{X_{..}} \right) & \cdots & \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{IJ}}{X_{I.}} - \frac{X_{.J}}{X_{..}} \right) \end{bmatrix}$$

- Por outro lado, temos:

$$\mathbf{D}_c = \begin{bmatrix} P_{.1} & 0 & \dots & 0 \\ 0 & P_{.2} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & P_{.J} \end{bmatrix} = \begin{bmatrix} \frac{X_{.1}}{X_{..}} & 0 & \dots & 0 \\ 0 & \frac{X_{.2}}{X_{..}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \frac{X_{.J}}{X_{..}} \end{bmatrix}$$

■ Assim:

$$\begin{aligned}
 W &= D_r^{1/2} (R - \mathbf{1}P'_c) D_c^{-1/2} = \\
 &= \begin{bmatrix} \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{11}}{X_{1.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{12}}{X_{1.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{1.}}{X_{..}}} \left(\frac{X_{1J}}{X_{1.}} - \frac{X_{.J}}{X_{..}} \right) \\ \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{21}}{X_{2.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{22}}{X_{2.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{2.}}{X_{..}}} \left(\frac{X_{2J}}{X_{2.}} - \frac{X_{.J}}{X_{..}} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{I1}}{X_{I.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{I2}}{X_{I.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{I.}}{X_{..}}} \left(\frac{X_{IJ}}{X_{I.}} - \frac{X_{.J}}{X_{..}} \right) \end{bmatrix} \\
 &\times \begin{bmatrix} \sqrt{\frac{X_{..}}{X_{.1}}} & 0 & \dots & 0 \\ 0 & \sqrt{\frac{X_{..}}{X_{.2}}} & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sqrt{\frac{X_{..}}{X_{.J}}} \end{bmatrix}
 \end{aligned}$$

■ Assim:

$$\begin{aligned}
 W &= \mathbf{D}_r^{1/2} (\mathbf{R} - \mathbf{1}\mathbf{P}'_c) \mathbf{D}_c^{-1/2} = \\
 &= \begin{bmatrix} \sqrt{\frac{X_{1.}}{X_{.1}}} \left(\frac{X_{11}}{X_{1.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{1.}}{X_{.2}}} \left(\frac{X_{12}}{X_{1.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{1.}}{X_{.j}}} \left(\frac{X_{1j}}{X_{1.}} - \frac{X_{.j}}{X_{..}} \right) \\ \sqrt{\frac{X_{2.}}{X_{.1}}} \left(\frac{X_{21}}{X_{2.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{2.}}{X_{.2}}} \left(\frac{X_{22}}{X_{2.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{2.}}{X_{.j}}} \left(\frac{X_{2j}}{X_{2.}} - \frac{X_{.j}}{X_{..}} \right) \\ \vdots & \vdots & \ddots & \vdots \\ \sqrt{\frac{X_{l.}}{X_{.1}}} \left(\frac{X_{l1}}{X_{l.}} - \frac{X_{.1}}{X_{..}} \right) & \sqrt{\frac{X_{l.}}{X_{.2}}} \left(\frac{X_{l2}}{X_{l.}} - \frac{X_{.2}}{X_{..}} \right) & \dots & \sqrt{\frac{X_{l.}}{X_{.j}}} \left(\frac{X_{lj}}{X_{l.}} - \frac{X_{.j}}{X_{..}} \right) \end{bmatrix}
 \end{aligned}$$

- Portanto:

$$W =$$

$$\begin{bmatrix} \left[X_{11} - \left(\frac{X_{1.} X_{.1}}{X_{..}} \right) \right] \frac{1}{(X_{1.} X_{.1})^{1/2}} & \dots & \left[X_{1J} - \left(\frac{X_{1.} X_{.J}}{X_{..}} \right) \right] \frac{1}{(X_{1.} X_{.J})^{1/2}} \\ \left[X_{21} - \left(\frac{X_{2.} X_{.1}}{X_{..}} \right) \right] \frac{1}{(X_{2.} X_{.1})^{1/2}} & \dots & \left[X_{2J} - \left(\frac{X_{2.} X_{.J}}{X_{..}} \right) \right] \frac{1}{(X_{2.} X_{.J})^{1/2}} \\ \vdots & \ddots & \vdots \\ \left[X_{I1} - \left(\frac{X_{I.} X_{.1}}{X_{..}} \right) \right] \frac{1}{(X_{I.} X_{.1})^{1/2}} & \dots & \left[X_{IJ} - \left(\frac{X_{I.} X_{.J}}{X_{..}} \right) \right] \frac{1}{(X_{I.} X_{.J})^{1/2}} \end{bmatrix}$$

- Dessa forma:

$$W_{ij} = \left[X_{ij} - \frac{X_{i.}X_{.j}}{X_{..}} \right] \frac{1}{(X_{i.}X_{.j})^{1/2}}$$

- Comparando o resultado da matriz \mathbf{W} com a estatística Q (teste de qui-quadrado), pode-se perceber que o quadrado dos elementos de \mathbf{W} são proporcionais ($Q = \sum_{i=1}^I \sum_{j=1}^J \left(\left[X_{ij} - \frac{X_{i.}X_{.j}}{X_{..}} \right] \frac{X_{..}}{(X_{i.}X_{.j})^{1/2}} \right)^2$) à contribuição de cada célula para a estatística de qui-quadrado.
- Como a estatística de qui-quadrado mede a associação entre as variáveis, então os elementos de \mathbf{W} fornecem também uma medida de associação para cada célula da tabela de contingência.
- Pode-se provar que $\mathbf{Z} = \mathbf{W}'$.

- Portanto, as matrizes $\mathbf{W}_{(I \times J)}$ e $\mathbf{Z}_{(J \times I)}$ dependem, respectivamente, dos perfis das linhas e dos perfis das colunas e também estão associadas à estrutura de dependência das variáveis presentes na tabela de contingência.
- O objetivo agora é tentar aproximar as duas matrizes através de um número menor de fatores, que possam representar apropriadamente cada uma delas, ou seja, por exemplo:

$$\mathbf{W}_{(I \times 2)}^* \approx \mathbf{W}; \mathbf{Z}_{(J \times 2)}^* \approx \mathbf{Z}$$

- (Decomposição do valor singular). Seja $\mathbf{A}_{(I \times J)}$, então podemos escrever :

$$\mathbf{A} = \mathbf{U}_{(I \times I)} \mathbf{\Lambda}_{(I \times J)} \mathbf{V}'_{(J \times J)}$$

em que

\mathbf{U} : colunas formadas pelos autovetores (ortonormalizados) de $\mathbf{A}\mathbf{A}'$.

\mathbf{V} : coluna formadas pelos autovetores (ortonormalizados) de $\mathbf{A}'\mathbf{A}$.

$\mathbf{\Lambda}$: matriz diagonal com sua diagonal dada por

$(\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{\min(I,J)})'$ (valores singulares), que são as raízes quadradas positivas dos autovalores maiores que zero, obtidos a partir das matrizes $\mathbf{A}\mathbf{A}'$ ou $\mathbf{A}'\mathbf{A}$ (os autovalores maiores que zero são iguais).

- Decomposição do valor singular de \mathbf{W} :

$$\mathbf{W} = \mathbf{U}_{W(I \times I)} \mathbf{\Lambda}_{W(I \times J)} \mathbf{V}'_{W(J \times J)}$$

- \mathbf{U}_W autovetores de $\mathbf{W}\mathbf{W}'$ e \mathbf{V}_W autovetores de $\mathbf{W}'\mathbf{W}$. $\mathbf{\Lambda}_W$ valores singulares de $\mathbf{W}\mathbf{W}'$ ou $\mathbf{W}'\mathbf{W}$, $\lambda_i \geq 0$.
- $(\mathbf{R} - \mathbf{1}\mathbf{P}'_c) = \mathbf{D}_r^{-1/2} \mathbf{U}_W \mathbf{\Lambda}_W \mathbf{V}'_W \mathbf{D}_c^{1/2}$.
- Resultado: A matriz \mathbf{A}^* que minimiza $\|\mathbf{A}^* - (\mathbf{R} - \mathbf{1}\mathbf{P}'_c)\|$ é dada por $\mathbf{A}^* = \mathbf{D}_r^{-1/2} \mathbf{U}_W^{(2)} \mathbf{\Lambda}_W^{(2)} \mathbf{V}_W'^{(2)} \mathbf{D}_c^{1/2}$ (considerando-se dois autovalores e autovetores).
- Em que: $\mathbf{U}_W^{(2)}$, $\mathbf{V}_W^{(2)}$ são as matrizes \mathbf{U}_W e \mathbf{V}_W , respectivamente, considerando, somente, as duas primeiras colunas destas e $\mathbf{\Lambda}_W^{(2)} = \text{diag}(\lambda_1, \lambda_2)$.

- Adicionalmente, pode-se provar que $\mathbf{F} = \mathbf{D}_{r(I \times I)}^{-1/2} \mathbf{U}_{W(I \times I)} \mathbf{\Lambda}_{W(I \times J)}$ definem coordenadas relacionadas aos perfis das linhas.
- Em outras palavras, $\mathbf{V}'_W \mathbf{D}_c^{1/2}$ correspondem à um tipo de “fator” (ou componente principal), enquanto que $\mathbf{D}_{r(I \times I)}^{-1/2} \mathbf{U}_{W(I \times I)} \mathbf{\Lambda}_{W(I \times J)}$ representam os coeficientes que associam os perfis das linhas à esses fatores.
- Como $\mathbf{Z} = \mathbf{W}'$, de forma análoga, temos que $(\mathbf{C} - \mathbf{1P}'_r) = \mathbf{D}_c^{-1/2} \mathbf{U}_Z \mathbf{\Lambda}_Z \mathbf{V}'_Z \mathbf{D}_r^{1/2}$ (exercício), em que $\mathbf{Z} = \mathbf{U}_{Z(J \times J)} \mathbf{\Lambda}_{Z(J \times I)} \mathbf{V}'_{Z(I \times I)}$ é a decomposição do valor singular de \mathbf{Z} .

- Adicionalmente, pode-se provar que $\mathbf{G} = \mathbf{D}_{c(J \times J)}^{-1/2} \mathbf{U}_{Z(J \times J)} \mathbf{\Lambda}_{Z(J \times I)}$ definem coordenadas relacionadas aos perfis das colunas.
- Ou seja, $\mathbf{V}'_Z \mathbf{D}_r^{1/2}$ correspondem à um tipo de “fator” (ou componente principal), enquanto que $\mathbf{D}_{c(J \times J)}^{-1/2} \mathbf{U}_{Z(J \times J)} \mathbf{\Lambda}_{Z(J \times I)}$ representam os coeficientes que associam os perfis das colunas à esses fatores.

- Calcular e dispor num gráfico de dispersão as seguintes quantidades (que correspondem à melhor representação bi dimensional de cada um dos conjunto de pontos associados aos perfis das linhas e aos perfis das colunas):

$$\mathbf{F}^{(2)} = \mathbf{D}_{r(I \times I)}^{-1/2} \mathbf{U}_{W(I \times 2)}^{(2)} \mathbf{\Lambda}_{W(2 \times 2)}^{(2)}$$

$$\mathbf{G}^{(2)} = \mathbf{D}_{c(J \times J)}^{-1/2} \mathbf{U}_{Z(J \times 2)}^{(2)} \mathbf{\Lambda}_{Z(2 \times 2)}^{(2)}$$

- Em que: $\mathbf{U}_Z^{(2)}$, $\mathbf{V}_Z^{(2)}$ são as matrizes \mathbf{U}_Z e \mathbf{V}_Z , respectivamente, considerando, somente, as duas primeiras colunas destas e $\mathbf{\Lambda}_Z^{(2)} = \text{diag}(\lambda_1, \lambda_2)$.

- Note, no entanto, que:

$$(R - \mathbf{1}P'_c) = D_r^{-1/2} U_W \Lambda_W V'_W D_c^{1/2}$$

$$D_r (R - \mathbf{1}P'_c) = D_r^{1/2} U_W \Lambda_W V'_W D_c^{1/2}$$

$$D_r R - D_r \mathbf{1}P'_c = D_r^{1/2} U_W \Lambda_W V'_W D_c^{1/2}$$

$$P - P_r P'_c = D_r^{1/2} U_W \Lambda_W V'_W D_c^{1/2}$$

- $P - P_r P'_c$: diferença entre probabilidades observadas e esperadas sob independência.

- Analogamente, note que:

$$(C - 1P'_r) = D_c^{-1/2} U_Z \Lambda_Z V'_Z D_r^{1/2}$$

$$D_c (C - 1P'_r) = D_c^{1/2} U_Z \Lambda_Z V'_Z D_r^{1/2}$$

$$D_c C - D_c 1P'_r = D_c^{1/2} U_Z \Lambda_Z V'_Z D_r^{1/2}$$

$$P - P_c P'_r = D_r^{1/2} U_Z \Lambda_Z V'_Z D_c^{1/2}$$

- $P - P_c P'_r$: diferença entre probabilidades observadas e esperadas sob independência.
- Em suma, estamos aproximando as “matrizes de distâncias” por suas decomposições de valores singulares bidimensionais.

- Variabilidade explicada.
- Inércia

$$\begin{aligned} & \text{tr} \left[\mathbf{D}_r^{-1/2} (\mathbf{P} - \mathbf{P}_r \mathbf{P}'_c) \mathbf{D}_c^{-1/2} \left(\mathbf{D}_r^{-1/2} (\mathbf{P} - \mathbf{P}_r \mathbf{P}'_c) \mathbf{D}_c^{-1/2} \right)' \right] \\ &= Q/n = \sum_{k=1}^{\min(I,J)} \lambda_k^2 \end{aligned}$$

em que $(\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_{\min(I,J)})$ são os valores singulares obtidos a partir da decomposição do valor singular de

$$\mathbf{D}_r^{-1/2} (\mathbf{P} - \mathbf{P}_r \mathbf{P}'_c) \mathbf{D}_c^{-1/2}.$$

- Para o caso bidimensional, a inércia será dada por: $\lambda_1^2 + \lambda_2^2$ e, a proporção de variabilidade explicada por $\frac{\lambda_1^2 + \lambda_2^2}{\sum_{k=1}^{\min(I,J)} \lambda_k^2}$.

- Análise de correspondência via R.
- Função *corresp* (pacote MASS) e *ca* (pacote ca).
- Função *ca*:
 - `resultCA <- ca(m.X)`
 - `inercia <- summary(resultCA)$scree`
 - `resultFCA <- plot(resultCA,xlab="componente 1",ylab="componente 2")`
 - `biplot(resultFCA$rows,resultFCA$cols,var.axes=FALSE,xlab="componente 1",ylab="componente 2",cex=1.2)`

■ Inércia

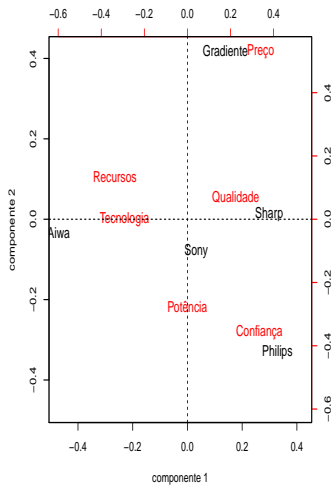
Valor singular	Inércia Princ.	Percentual	Percentual acum.
0,2678	0,0717	52,71	52,71
0,2228	0,0497	36,49	89,20
0,1023	0,0105	7,69	96,88
0,0651	0,0042	3,12	100,00
Total	0,1368	100,00	-

- Vemos que duas componentes explicam aproximadamente 89,20% da variabilidade dos dados.

■ Componentes:

	Categoria	Componente 1	Componente 2
Perfil das Linhas (F)	Sony	0,032	-0,0804
	Aiwa	-0,4703	-0,0326
	Gradiente	0,1386	0,4188
	Philips	0,3286	-0,3301
	Sharp	0,2989	0,0131
Perfil das Colunas (G)	Qualidade	0,2246	0,0713
	Tecnologia	-0,2924	0,0006
	Potência	-0,002	-0,2773
	Recursos	-0,3368	0,135
	Preço	0,3399	0,532
	Confiança	0,3342	-0,3578

■ Biplot:



- O gráfico mostra que a marca Gradiente está mais relacionada com os preços, ou seja, consumidores da marca Gradiente optam por ela devido ao preço.
- A marca Sharp está bastante relacionada com qualidade do som e a marca Philips está relacionada com a confiança na marca.
- A Sony está mais relacionada com Potência do som, qualidade do som e tecnologia.
- Por fim, a Aiwa está bem relacionada com a tecnologia avançada e recursos técnicos.

Exemplo 7: espécies extintas

- Diz respeito ao número de espécies extintas desde o ano de 1600 aproximadamente.
- Os dados correspondem ao número de certas espécies extintas por continente (excluindo a Antártica).

Cont.	Espécie						Total
	Moluscos	Insetos	Peixes	Répteis	Aves	Mamíferos	
Ásia	0	1	0	0	14	5	20
Europa	2	1	0	0	4	3	10
América	40	9	31	9	18	30	137
Oceania	81	47	1	2	51	23	205
África	68	3	0	12	39	5	127
Total	191	61	32	23	126	66	499

OBS: Cont. - continente. Estatística de qui-quadrado: $Q = 192,51$,
 $p\text{-valor} < 0,0001$.

■ Perfis das linhas ($\times 100$)

Cont.	Espécie						Total
	Moluscos	Insetos	Peixes	Répteis	Aves	Mamíferos	
Ásia	0,00	5,00	0,00	0,00	70,00	25,00	100,00
Europa	20,00	10,00	0,00	0,00	40,00	30,00	100,00
América	29,20	6,57	22,63	6,57	13,14	21,90	100,00
Oceania	39,51	22,93	0,49	0,98	24,88	11,22	100,00
África	53,54	2,36	0,00	9,45	30,71	3,94	100,00
Total	38,28	12,22	6,41	4,61	25,25	13,23	100,00

■ Perfis das colunas ($\times 100$)

Cont.	Espécie						Total
	Moluscos	Insetos	Peixes	Répteis	Aves	Mamíferos	
Ásia	0,00	1,64	0,00	0,00	11,11	7,58	4,01
Europa	1,05	1,64	0,00	0,00	3,17	4,55	2,00
América	20,94	14,75	96,88	39,13	14,29	45,45	27,45
Oceania	42,41	77,05	3,12	8,70	40,48	34,85	41,08
África	35,60	4,92	0,00	52,17	30,95	7,58	25,45
Total	100,00	100,00	100,00	100,00	100,00	100,00	100,00

■ Inércia

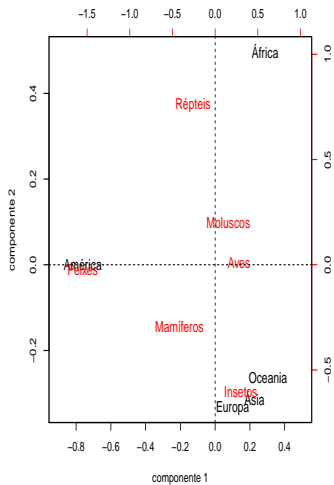
Valor singular	Inércia Princ.	Percentual	Percentual acum.
0,4698	0,2207	57,20	57,20
0,3130	0,0980	25,40	82,60
0,2553	0,0652	16,89	99,49
0,0443	0,0020	0,51	100,00

- Vemos que duas componentes explicam aproximadamente 82,60% da variabilidade dos dados.

■ Componentes:

	Categoria	Componente 1	Componente 2
Perfil das Linhas (F)	Ásia	0,2264	-0,3105
	Europa	0,1009	-0,3348
	América	-0,7620	0,0009
	Oceania	0,3039	-0,2632
	África	0,2877	0,4991
Perfil das Colunas (G)	Moluscos	0,1550	0,2005
	Insetos	0,3007	-0,6028
	Peixes	-1,5511	-0,0234
	Répteis	-0,2588	0,7600
	Aves	0,2801	0,0095
	Mamíferos	-0,4191	-0,2947

■ Biplot:



- O biplot mostra que a Oceania, Ásia e a Europa estão bastante próximas e mais relacionadas com a extinção de insetos e, de uma forma menos intensa, com a extinção de mamíferos.
- Na África houve uma extinção em maior número de espécies de moluscos e répteis e, com uma menor intensidade, de aves.
- No continente Americano, a extinção espécies de peixes e mamíferos foi mais acentuada do que as demais.