

# Análise Multivariada

Lupércio França Bessegato  
Dep. Estatística/UFJF

## Introdução

### Roteiro

1. Introdução
2. Distribuições de Probabilidade Multivariadas
3. Representação de Dados Multivariados
4. Testes de Significância c/ Dados Multivariados
5. Análise de Componentes Principais
6. Análise Fatorial
7. Análise de Agrupamentos
8. Análise de Correlação Canônica
9. Referências

Análise Multivariada - 2022

2

### Classificação e Agrupamento

- Classificar:
  - ✓ Número de grupos é conhecido e o objetivo é alocar novas observações a um desses grupos
- Agrupar:
  - ✓ Não há suposições sobre o número de grupos ou sobre a estrutura dos grupos
    - Técnica mais primitiva

Análise Multivariada - 2022

4

### **Análise de Agrupamento e Análise Discriminante**

- Análise de Agrupamentos
  - √ Dividir os elementos da amostra (ou população) em grupos, de maneira que:
    - Elementos de um grupo são similares entre si
    - Elementos de grupos diferentes sejam heterogêneos em relação a essas características

Análise Multivariada - 2022

5

- Análise discriminante:
  - √ Classificação de elementos de amostra (população)
    - Grupos são pré-definidos
  - √ Procedimento:
    - Regra de classificação

Análise Multivariada - 2022

6

### **Análise de Agrupamentos**

### **Análise de Agrupamentos**

- Procurar por uma estrutura de grupos “naturais” dos dados
  - √ É uma importante técnica exploratória
- Objetivo básico:
  - √ Descobrir agrupamentos naturais dos itens (ou variáveis)
- Em geral, somos capazes de agrupar visualmente objetos em gráficos

Análise Multivariada - 2022

8

- São necessários:
  - √ Medidas de similaridade (ou distância)
  - √ Desenvolvimento de escala quantitativa para medir associação (similaridade) entre os dados
  - √ Algoritmos para ordenar objetos em grupos

Análise Multivariada - 2022

9

## Medidas de Similaridade

- Há muita subjetividade na escolha de uma medida de similaridade
- Considerações importantes:
  - √ Natureza das variáveis
    - (discreta, contínua, binária)
  - √ Escala das medidas
    - (nominal, ordinal, intervalar, razão)

Análise Multivariada - 2022

10

- Agrupamentos de itens (unidades ou casos)
  - √ Proximidade é usualmente indicada por algum tipo de distância
- Agrupamento de variáveis:
  - √ Usualmente são agrupadas com base em coeficientes de correlação ou medidas de associação

Análise Multivariada - 2022

11

## Distâncias para Pares de Itens

- Sejam as observações:  
 $\sqrt{\mathbf{x}'} = [x_1, x_2, \dots, x_p]$  e  $\mathbf{y}' = [y_1, y_2, \dots, y_p]$
- Distância Euclidiana:

$$\begin{aligned}d(\mathbf{x}, \mathbf{y}) &= \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 \dots (x_p - y_p)^2} \\ &= \sqrt{(\mathbf{x} - \mathbf{y})'(\mathbf{x} - \mathbf{y})}\end{aligned}$$

Análise Multivariada - 2022

12

- Distância generalizada ou ponderada:

$$d(\mathbf{x}, \mathbf{y}) = \sqrt{(\mathbf{x} - \mathbf{y})' \mathbf{A} (\mathbf{x} - \mathbf{y})}$$

- ✓  $\mathbf{A}$  é matriz de ponderação positiva definida
- ✓  $\mathbf{A} = \mathbf{S}^{-1}$  (distância de Mahalanobis)
  - Não podem ser calculadas sem conhecimento prévio dos grupos
- ✓ Se  $\mathbf{A} = \mathbf{I}$  (distância Euclidiana)
- ✓ Se  $\mathbf{A} = \text{diagonal}(1/p)$  (distância Euclidiana média)

Análise Multivariada - 2022

13

- Métrica de Minkowski:

$$d(\mathbf{x}, \mathbf{y}) = \left[ \sum_{i=1}^p w_i |x_i - y_i|^m \right]^{1/m}$$

- ✓  $w_i$ : peso de ponderação para as variáveis
- ✓  $m = 1$ ,  $d(\mathbf{x}, \mathbf{y})$  mede distância “city block” ou Manhattan
- ✓  $m = 2$ ,  $d(\mathbf{x}, \mathbf{y})$  é a distância Euclidiana
- ✓ variar  $m$  muda a ponderação dada a diferenças maiores ou menores
- ✓ A métrica de Minkowski é menos afetada pela presença de valores discrepantes na amostra do que a distância Euclidiana.

Análise Multivariada - 2022

14

### Métricas para Variáveis Não-Negativas

- Métrica de Canberra:

$$d(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^p \frac{|x_i - y_i|}{x_i + y_i}$$

- Métrica de Czekanowski:

$$d(\mathbf{x}, \mathbf{y}) = 1 - \frac{2 \sum_{i=1}^p \min(x_i, y_i)}{\sum_{i=1}^p x_i + y_i}$$

Análise Multivariada - 2022

15

### Distância

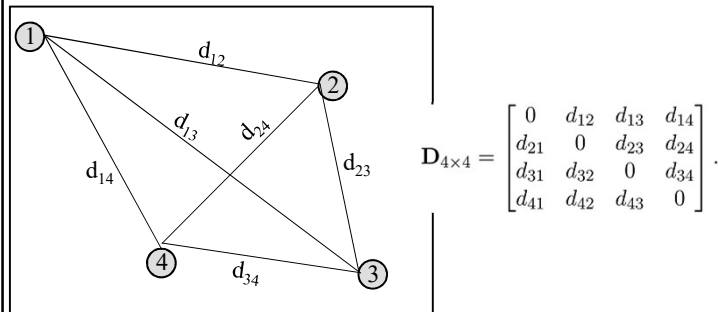
- Qualquer medida de distância  $d(\mathbf{P}, \mathbf{Q})$  entre dois pontos  $\mathbf{P}$  e  $\mathbf{Q}$  é válida, desde que satisfaça as seguintes propriedades.  $\mathbf{R}$  é um ponto intermediário:
  - ✓  $d(\mathbf{P}, \mathbf{Q}) = d(\mathbf{Q}, \mathbf{P})$
  - ✓  $d(\mathbf{P}, \mathbf{Q}) > 0$  se  $\mathbf{P} \neq \mathbf{Q}$
  - ✓  $d(\mathbf{P}, \mathbf{Q}) = 0$  se  $\mathbf{P} = \mathbf{Q}$
  - ✓  $d(\mathbf{P}, \mathbf{Q}) \leq d(\mathbf{P}, \mathbf{R}) + d(\mathbf{R}, \mathbf{Q})$  – desigualdade triangular

Análise Multivariada - 2022

16

### Exemplo – Matriz de Distâncias

- Esquema de armazenamento de distâncias



Análise Multivariada - 2022

17

### Exemplo

- Renda mensal (em SM) e idade

Indivíduo	Renda	Idade
A	9,60	28
B	8,40	31
C	2,40	42
D	18,20	38
E	3,90	25
F	6,40	41
Média	8,15	34,17
Desvio padrão	5,61	7,14

$$\bar{\mathbf{x}} = \begin{bmatrix} 8,15 \\ 34,17 \end{bmatrix}$$

$$\mathbf{S}_{2 \times 2} = \begin{bmatrix} 5,61^2 & 2,15 \\ 2,15 & 7,14^2 \end{bmatrix}$$

Análise Multivariada - 2022

18

- Distância entre os indivíduos A e B:

√ Distância euclidiana:

$$d(\mathbf{X}_A, \mathbf{X}_B) = \sqrt{(9,60 - 8,40)^2 + (28 - 31)^2} = 3,23.$$

√ Distância de Mahalanobis:

$$d(\mathbf{X}_A, \mathbf{X}_B) = \left( [1, 2 \quad -3] \mathbf{S}^{-1} \begin{bmatrix} 1,2 \\ -3 \end{bmatrix} \right)^{1/2}$$

$$\mathbf{S}^{-1} = \begin{bmatrix} 31,47 & 2,15 \\ 2,15 & 50,97 \end{bmatrix}^{-1} = \begin{bmatrix} 0,032 & -0,0013 \\ -0,0013 & 0,019 \end{bmatrix}$$

$$d(\mathbf{X}_A, \mathbf{X}_B) = \sqrt{(1,2)^2(0,032) - 2(1,2)(-3)(-0,0013) + (-3)^2(0,019)} = 0,46$$

Análise Multivariada - 2022

19

√ Distância euclidiana média:

$$d(\mathbf{X}_A, \mathbf{X}_B) = \left( [1, 2 \quad -3] \mathbf{A}^{-1} \begin{bmatrix} 1,2 \\ -3 \end{bmatrix} \right)^{1/2}$$

$$\mathbf{A}^{-1} = \begin{bmatrix} 31,47 & 0 \\ 0 & 50,97 \end{bmatrix}^{-1} = \begin{bmatrix} 0,032 & 0 \\ 0 & 0,019 \end{bmatrix}$$

$$d(\mathbf{X}_A, \mathbf{X}_B) = \sqrt{(1,2)^2(0,032) + (-3)^2(0,019)} = 0,47$$

√ Distância de Minkowsky:

$$- w_i = 1 \text{ e } m = 1$$

$$d(\mathbf{X}_A, \mathbf{X}_B) = \left[ \sum_{i=A}^B w_i |x_i - y_i|^m \right]^{1/m}$$

$$d(\mathbf{X}_A, \mathbf{X}_B) = |9,60 - 8,40| + |28 - 31| = 4,20.$$

Análise Multivariada - 2022

20

- Itens representados por medidas qualitativas

✓ os pares de itens são frequentemente comparados com base na presença ou ausência de certas características

✓ Itens similares têm mais características comuns que os itens dissimilares

✓ Presença ou ausência de característica é descrita por variável indicadora (binária):

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>
Item i	1	0	0	1	1
Item j	1	1	0	1	0

Análise Multivariada - 2022

21

- Para  $j = 1, 2, \dots, p$ , sejam:

$x_{ij}$ : escore da  $j$ -ésima variável do  $i$ -ésimo item

$x_{kj}$ : escore da  $j$ -ésima variável do  $k$ -ésimo item

$$(x_{ij} - x_{kj})^2 = \begin{cases} 0 & \text{se } x_{ij} = x_{kj} = 1 \text{ ou } x_{ij} = x_{kj} = 0 \\ 1 & \text{se } x_{ij} \neq x_{kj} \end{cases}$$

✓ A distância Euclidiana  $\sum_{i=1}^p (x_{ij} - x_{kj})^2$  é a contagem das discordâncias

✓ Grandes distâncias correspondem a muitas discordâncias

✓ Essa medida de similaridade pondera igualmente concordâncias e discordâncias

Análise Multivariada - 2022

22

- No exemplo:

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>
Item i	1	0	0	1	1
Item j	1	1	0	1	0

$$\sum_{i=1}^p (x_{ij} - x_{kj})^2 = (1-1)^2 + (0-1)^2 + (0-0)^2 + (1-1)^2 + (1-0)^2 = 2$$

• Muitas vezes uma concordância 1–1 é uma indicação mais forte de similaridade que uma concordância 0–0

Análise Multivariada - 2022

23

## Coefficientes de Similaridade

- Comparação variáveis qualitativas

✓ Presença ou ausência de característica (0 ou 1)

✓ Elementos comuns devem ter em comum mais itens similares que dissimilares

- Coeficientes de Similaridade

✓ Coeficiente de concordância simples

✓ Coeficiente de concordância positiva

✓ Coeficiente de concordância de Jaccard

✓ Distância euclidiana média

✓ Outros

Análise Multivariada - 2022

24

• Coeficiente de Concordância Simples:

		(k)			
		1	0		
(i)	1	a	b	a + b	$s(i, k) = \frac{a + d}{p}$
	0	c	d	c + d	
		a + c	b + d	p	

- ✓ Pondera da mesma maneira as concordâncias
- ✓ Quanto maior o valor de  $s(i, k)$  maior a similaridade

Análise Multivariada - 2022

25

• Coeficiente de Concordância Positiva:

		(k)			
		1	0		
(i)	1	a	b	a + b	$s(i, k) = \frac{a}{p}$
	0	c	d	c + d	
		a + c	b + d	p	

- ✓ Considera apenas os empates 1-1
  - Empates 0-0 não representam necessariamente uma concordância
- ✓ Quanto maior o valor de  $s(i, k)$  maior a similaridade

Análise Multivariada - 2022

26

• Coeficiente de Concordância de Jaccard:

- ✓ Usado em ecologia (Ludwig; reynolds, 1988)

		(k)			
		1	0		
(i)	1	a	b	a + b	$s(i, k) = \frac{a}{a + b + c}$
	0	c	d	c + d	
		a + c	b + d	p	

- ✓ Semelhante ao coeficiente de concordância positiva
  - Empates 0-0 tratados como irrelevantes

Análise Multivariada - 2022

27

• Outros Coeficientes de Concordância:

- ✓ Sokol e Sneath, 1963:  $s(i, k) = \frac{2(a + d)}{2(a + d) + b + c}$

- O dobro da ponderação para as concordâncias

- ✓ Roger e Tanimoto, 1960:  $s(i, k) = \frac{a + d}{a + d + 2(b + c)}$

- O dobro da ponderação para as discordâncias

- ✓ Czeakanowsky/Sorensen-Dice:  $s(i, k) = \frac{2a}{2a + b + c}$

- Pondera o dobro das concordâncias, sem considerar as discordâncias

Análise Multivariada - 2022

28

✓ Andersberg, 1973:

$$s(i, k) = \frac{a}{a + 2(b + c)}$$

- Pondera o dobro das discordâncias
- Não considera os empates 0-0

✓ Kulczynski I, 1927:

$$s(i, k) = \frac{a}{b + c}$$

- Razão das concordâncias com as discordâncias
- Não considera os empates 0-0

Análise Multivariada - 2022

29

## • Distância Euclidiana Média

$$d(i, k) = \left[ \frac{1}{p} \sum_{j=1}^p (x_{ij} - x_{kj})^2 \right]^{1/2} = \frac{\# \text{ pares discordantes}}{\# \text{ pares}}$$

✓ Quanto menor o valor da distância, maior a similaridade

✓ É um índice de discordância ou dissimilaridade

Análise Multivariada - 2022

30

## • Conjunto de Animais (Andreas, 2003)

✓ 16 animais descritos por características arbitrárias

	Feature	Dove	Hen	Duck	Goose	Owl	Hawk	Eagle	Fox	Dog	Wolf	Cat	Tiger	Lion	Horse	Zebra	Cow
Is	small	1	1	1	1	1	1	0	0	0	0	1	0	0	0	0	0
	medium	0	0	0	0	0	0	1	1	1	1	0	0	0	0	0	0
	big	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
has	2 legs	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
	4 legs	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
	hair	0	0	0	0	0	0	0	1	1	1	1	1	1	1	1	1
	hooves	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1
	mane	0	0	0	0	0	0	0	0	1	0	0	1	1	1	1	0
	feathers	1	1	1	1	1	1	1	0	0	0	0	0	0	0	0	0
likes to	hunt	0	0	0	0	1	1	1	1	0	1	1	1	1	0	0	0
	run	0	0	0	0	0	0	0	0	1	1	0	1	1	1	1	0
	fly	1	0	0	1	1	1	1	0	0	0	0	0	0	0	0	0
	swim	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0

Table 4 - Animals Data Set

Análise Multivariada - 2022

31

## • Algumas medidas de similaridade em pares

account d	Dove/Hen	Dove/Tiger	Tiger/Cow	Dove/Cow
Simple Matching	0,92	0,31	0,77	0,38
Russel and Rao	0,23	0,00	0,23	0,00
Rogers and Tanimoto	0,86	0,18	0,63	0,24
Hamman	0,85	-0,38	0,54	-0,23
Ochiai II	0,47	0,00	0,39	0,00
Sokal	0,96	0,47	0,87	0,56
Don't account d				
Jaccard	0,75	0,00	0,50	0,00
Anderberg	0,60	0,00	0,33	0,00
Kulczynski I	3,00	0,00	1,00	0,00
Kulczynski II	0,88	0,00	0,68	0,00
Sorensen-Dice	0,86	0,00	0,67	0,00
Ochiai	0,87	0,00	0,67	0,00

Table 5 - Some similarity measures values

✓ 2 grupos: aves e mamíferos

✓ Aves, carnívoros e herbívoros

Análise Multivariada - 2022

32



## Monotonicidade

- Relacionamento monotônico de coeficientes
  - ✓ Mantém a ordenação relativa das similaridades
  - ✓ Alguns procedimentos de agrupamento não são afetados se a definição de similaridade é mudada
- Coeficientes relacionados monotonicamente:
  - ✓ Condordância simples, Sokol-Sneath e Roger-Tanomoto
  - ✓ Jaccard, Czeakanowsky-Sorensen-Dice e Andersberg

Análise Multivariada - 2022

33

## Exemplo 12.1

- Cálculo dos valores de coeficiente de similaridade

✓ Características de indivíduos

Indivíduo	Altura	Peso	Olhos	Cabelo	Mão	Sexo
1	173	65	Verdes	Loiro	Destro	F
2	185	82	Castanhos	Castanhos	Destro	M
3	170	75	Azuis	Loiro	Destro	M
4	163	54	Castanhos	Castanho	Destro	F
5	193	95	Castanhos	Castanho	Canhoto	M

X<sub>1</sub>: 1, altura ≥ 180 cm  
0, altura < 180 cm

X<sub>3</sub>: 1, olhos castanhos  
0, caso contrário

X<sub>5</sub>: 1, destro  
0, canhoto

X<sub>2</sub>: 1, peso ≥ 70 kg  
0, peso < 70 kg

X<sub>4</sub>: 1, cabelos loiros  
0, caso contrário

X<sub>6</sub>: 1, Feminino  
0, Masculino

Análise Multivariada - 2022

36

✓ Indivíduos 1 e 2:

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
1	0	0	0	1	1	1
2	1	1	1	0	1	0

✓ Coeficiente de concordância simples

$$s(1,2) = \frac{a+d}{p} = \frac{1+0}{6} = \frac{1}{6}$$

✓ Matriz de similaridades:

– Mais similares: 2 e 5

– Menos similares: 1 e 5

✓ Grupos: (1, 3, 4) e (2, 5)

$$S = \begin{bmatrix} 1 & 2 & 3 & 4 & 5 \\ 1 & \frac{1}{6} & 0 & 0 & 0 \\ 2 & \frac{1}{6} & 1 & 0 & 0 \\ 3 & 0 & \frac{1}{6} & 1 & 0 \\ 4 & 0 & \frac{1}{6} & 1 & 0 \\ 5 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Análise Multivariada - 2022

37

## Similaridades & Distâncias

- Sempre é possível construir similaridades a partir de distâncias

✓ Exemplo:  $\tilde{s}(i,k) = \frac{1}{1+d(i,k)}$ ,  $0 < \tilde{s}(i,k) \leq 1$   
– d<sub>ik</sub>: distância entre os itens i e k

- Nem sempre é possível construir medidas que satisfaçam as propriedades de distância

- Matriz de similaridades deve ser positiva definida
- Padronizar a máxima similaridade em 1 (s<sub>ii</sub> = 1)

✓ Nessas condições, d(i,k) tem propriedades de distância:

$$d(i,k) = \sqrt{2(1 - \tilde{s}(i,k))}$$

Análise Multivariada - 2022

38

- Há muitas maneiras de medir similaridade:
  - √ Pares de itens:
    - Distâncias ou coeficientes de similaridade
  - √ Pares de variáveis
    - Coeficientes de correlação
- Pode-se também agrupar através de frequência

Análise Multivariada - 2022

42

## Exemplo 12.2

- O significado das palavras muda ao longo da história
  - √ O significado dos números constitui uma exceção
- Uma primeira comparação de línguas poderia ser baseada nos numerais

Análise Multivariada - 2022

43

## • Numerais em 11 línguas

English (E)	Norwegian (N)	Danish (D)	Dutch (Du)	German (G)	French (F)	Spanish (S)	Italian (I)	Polish (P)	Hungarian (H)	Finnish (F)
one	en	en	een	eins	un	uno	uno	jeden	egy	yksi
two	to	to	twee	zwei	deux	dos	due	dwa	ketto	kaksi
three	tre	tre	drie	trois	tres	tre	tre	trzy	három	kolme
four	fire	fire	vier	vier	quatre	cuatro	quattro	cztery	negy	neua
five	fem	fem	vijf	funf	cinq	cinco	cinq	piec	öt	viisi
six	seks	seks	zes	sechs	six	seis	sei	szesc	hat	kuusi
seven	sju	syv	zeven	sieben	sept	siete	sette	siedem	het	seitsemän
eight	atte	otte	acht	acht	huit	ocho	otto	osiem	nyolc	kahdeksan
nine	ni	ni	neun	neuf	nueve	nove	dziesięć	kilenc	yhdeksän	
ten	ti	ti	tien	zehn	dix	diez	dieci	dziesięć	tíz	kymmenen

√ Comparação das línguas pela 1ª. letra dos números

- Números concordantes: tem a mesma 1ª. letra
- Números discordantes: caso contrário

Análise Multivariada - 2022

44

	E	N	D	Du	G	F	S	I	P	H	F
E	10										
N	8	10									
D	8	9	10								
Du	3	5	4	10							
G	4	6	5	5	10						
F	4	4	4	1	3	10					
S	4	4	5	1	3	8	10				
I	4	4	5	1	3	9	9	10			
P	3	3	4	0	2	5	7	6	10		
H	1	2	2	2	1	0	0	0	0	10	
F	1	1	1	1	1	1	1	1	1	2	10

√ Inglês e Norueguês – 1ª.s letras: 8 em 10

√ Inglês, norueguês, dinamarquês, holandês e alemão  
– aparentam formar um grupo

√ Francês, espanhol, italiano e polonês  
– podem ser agrupados

√ Húngaro e finlandês parecem estar sozinhos

Análise Multivariada - 2022

45

### Métodos de Agrupamentos Hierárquicos

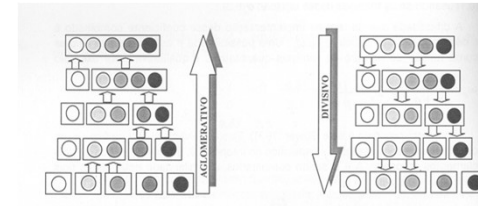
- Raramente podemos examinar todas as possibilidades de agrupamentos
  - ✓ Há algoritmos de agrupamento que não têm de verificar todas as configurações
- Técnicas de Agrupamento Hierárquicas
  - ✓ Procedimentos que realizam uma série de sucessivas fusões (ou uma série de sucessivas divisões)

Análise Multivariada - 2022

46

### • Técnicas Hierárquicas:

- ✓ Aglomerativas
- ✓ Divisivas



- ✓ Em geral, são usadas em análises exploratórias dos dados com o objetivo de:
  - identificar possíveis agrupamentos
  - estimar o valor provável do número de grupos  $g$

Análise Multivariada - 2022

47

### • Técnicas Não-Hierárquicas:

- ✓ É necessário que o valor do número de grupos já esteja pré-especificado pelo pesquisador

Análise Multivariada - 2022

49

### Métodos Hierárquicos Aglomerativos

1. Cada elemento constitui um cluster de tamanho 1
  - ✓ Há  $n$  clusters
2. Em cada estágio do algoritmo os pares de conglomerados mais similares são combinados (novo conglomerado)
  - ✓ Em cada estágio do processo, o número de conglomerados vai sendo diminuído

Análise Multivariada - 2022

50

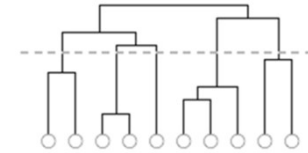
### 3. Propriedade de Hierarquia:

- ✓ Em cada estágio do algoritmo, cada novo conglomerado formado é um agrupamento de conglomerados formados nos estágios anteriores
  - Se 2 elementos aparecem juntos em algum estágio do processo, eles aparecerão juntos em todos os estágios subsequentes
- ✓ Uma vez unidos, estes elementos não poderão ser separados

Análise Multivariada - 2022

51

### 4. Dendrograma (ou Dendrograma):



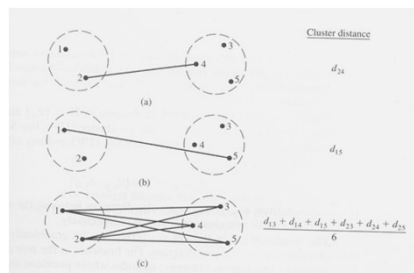
- ✓ Representa a árvore (ou história) do agrupamento
  - Escala Vertical: nível de similaridade (ou dissimilaridade)
  - Eixo Horizontal: elementos amostrais na ordem relacionada à história do agrupamento

Análise Multivariada - 2022

52

## Métodos de Agrupamentos

- Medida de similaridade (ou distância) entre 2 conglomerados

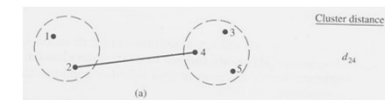


Análise Multivariada - 2022

54

### • Método de Ligação Simples (Single Linkage):

- ✓ Similaridade entre dois conglomerados é definida pelos dois elementos mais parecidos entre si
  - distância mínima ou vizinho mais próximo



$$C_1 = \{X_1, X_2\} \text{ e } C_2 = \{X_3, X_4, X_5\}$$

$$d(C_1, C_2) = \min\{d(X_j, X_k)\}, \quad j \neq k, \quad j = 1, 2 \text{ e } k = 3, 4, 5$$

- ✓ Em cada estágio do processo de agrupamento os dois conglomerados que são mais similares (mais próximos) são combinados em um único *cluster*.

Análise Multivariada - 2022

56

### Exemplo 12.3

- Matriz de Distâncias (**D**):

$$\sqrt{\min\{d_{ik}\}} = d(5,3) = 2$$

$$\sqrt{\text{cluster (35)}}$$

$$d(35, 1) = \min\{d(3, 1), d(5, 1)\} = \min\{3, 11\} = 3$$

$$d(35, 2) = \min\{d(3, 2), d(5, 2)\} = \min\{7, 10\} = 7$$

$$d(35, 4) = \min\{d(3, 4), d(5, 4)\} = \min\{9, 8\} = 8$$

$$\sqrt{\text{cluster (135)}}$$

$$d(135, 2) = \min\{d(35, 2), d(1, 2)\} = \min\{7, 9\} = 7$$

$$d(135, 4) = \min\{d(35, 4), d(1, 4)\} = \min\{8, 6\} = 6$$

$$\sqrt{\text{cluster (1354)}}$$

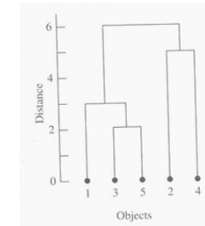
$$d(1354, 2) = d(135, 2) = 7$$

	1	2	3	4	5
1	0				
2	9	0			
3	3	7	0		
4	6	5	9	0	
5	11	10	2	8	0

Análise Multivariada - 2022

58

- Dendrograma:



Os resultados intermediários são o principal interesse

Análise Multivariada - 2022

59

### Exemplo 12.4

- Numerais em 11 línguas (continuação 12.2)

	E	N	D	Du	G	F	S	I	P	H	F
E	0										
N	2	0									
D	2	1	0								
Du	7	5	6	0							
G	6	4	5	5	0						
F	6	6	6	9	7	0					
S	6	6	5	9	7	2	0				
I	6	6	5	9	7	1	1	0			
P	7	7	6	10	8	5	3	4	0		
H	9	8	8	8	9	10	10	10	10	0	
F	9	9	9	9	9	9	9	9	9	8	0

Menores distâncias entre pares de línguas:

- $d(D,N)=1$ ;  $d(I,F)=1$ ;  $d(I,S)=1$
- Como  $d(F,S)=2$ , podemos fundir apenas IF ou IS

Análise Multivariada - 2022

61

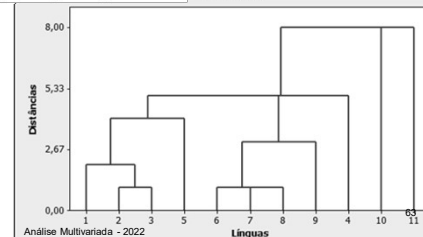
- Análise de Aglomerado – Ligação Simples

Cluster Analysis of Observations: MI

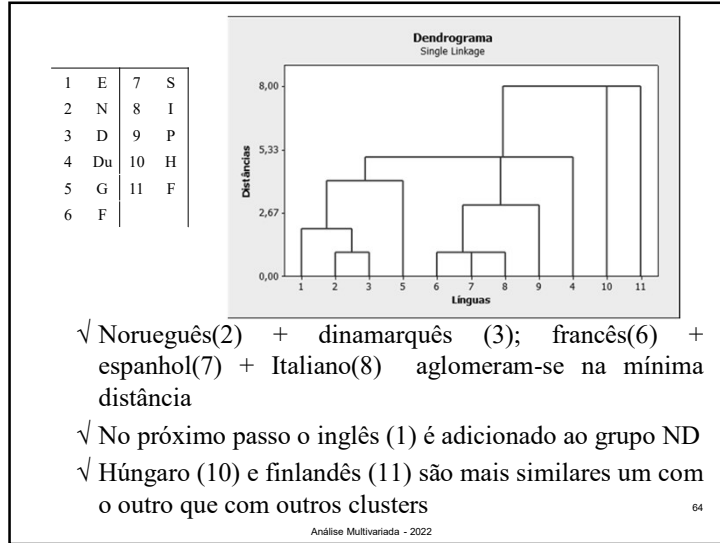
Single Linkage  
Agglomeration Steps

Step	Number of clusters	Similarity level	Distance level	Clusters joined	New cluster	Number of obs. in new cluster
1	10	90	1	7 8	7	2
2	9	90	1	6 7	6	3
3	8	90	1	2 3	2	2
4	7	80	2	1 2	1	3
5	6	70	3	6 9	6	4
6	5	60	4	1 5	1	4
7	4	50	5	1 6	1	8
8	3	50	5	1 4	1	9
9	2	20	8	10 11	10	2
10	1	20	8	1 10	1	11

Dendrograma  
Single Linkage

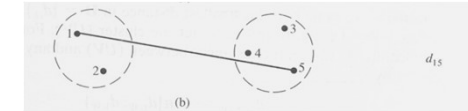


Análise Multivariada - 2022



• Método de Ligação Completa (Complete Linkage):

- √ Similaridade entre dois conglomerados é definida pelos dois elementos menos parecidos entre si
- distância máxima ou vizinho mais distante



$$C_1 = \{X_1, X_2\} \text{ e } C_2 = \{X_3, X_4, X_5\}$$

$$d(C_1, C_2) = \max\{d(X_j, X_k)\}, \quad j \neq k, \quad j = 1, 2 \text{ e } k = 3, 4, 5$$

- √ Em cada estágio, a distância (similaridade) entre os clusters é determinada pela distância (similaridade) entre os dois elementos, em cada cluster, que são mais distantes.

Análise Multivariada - 2022

- Garante que todos os itens em cada cluster estão com a máxima distância (mínima similaridade) entre eles.

Análise Multivariada - 2022

Exemplo 12.7

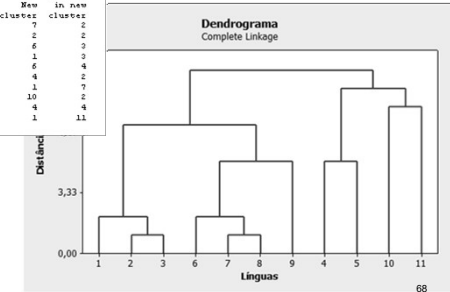
• Numerais em 11 línguas

√ Ligação completa

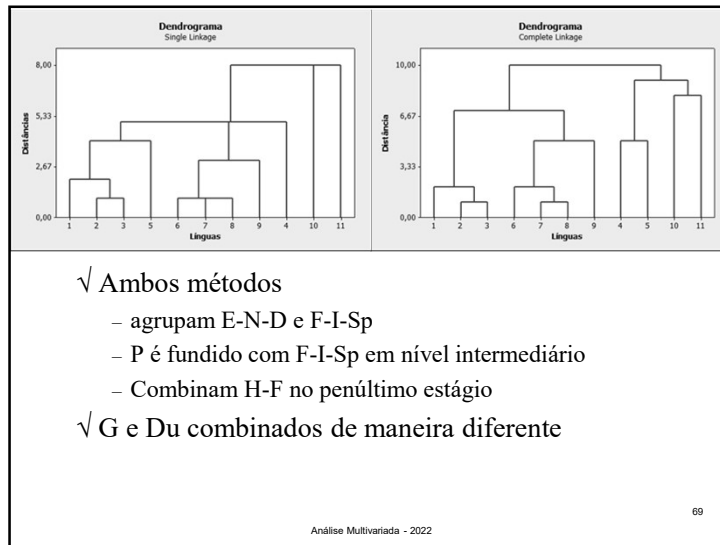
Cluster Analysis of Variables: M1

Complete Linkage  
Agglomeration Steps

Step	Number of clusters	Similarity level	Distance	Clusters joined	Number of obs. in new cluster
1	10	90	1	7 8	2
2	9	90	1	2 3	2
3	8	80	2	6 7	2
4	7	80	2	1 2	3
5	6	50	5	5 9	4
6	5	50	5	4 5	4
7	4	20	7	1 6	7
8	3	20	6	10 11	2
9	2	10	9	4 10	4
10	1	0	10	1 4	11



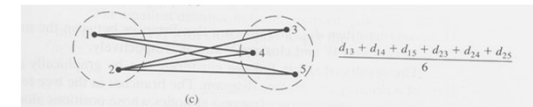
Análise Multivariada - 2022



• Método da Média das Distâncias (AverageLinkage):

√ Similaridade entre dois conglomerados é definida pela distância média de todos os pares de itens

– cada membro do par pertence a grupos diferentes



$C_1 = \{X_1, X_2\}$  e  $C_2 = \{X_3, X_4, X_5\}$

$$d(C_1, C_2) = \sum_{j \in C_1} \sum_{k \in C_2} \frac{n_1 n_2}{n_1 + n_2} d(X_j, X_k)$$

–  $n_1$ : quantidade de elementos do *cluster* 1

–  $n_2$ : quantidade de elementos do *cluster* 2.

Análise Multivariada - 2022

70

- Podem ser usadas distâncias ou similaridades
- Pode ser usado para agrupar variáveis e itens
- Mudanças na atribuição de distâncias (similaridade) podem afetar o arranjo da configuração final de *clusters*, mesmo que as alterações preservem as ordenações relativas.

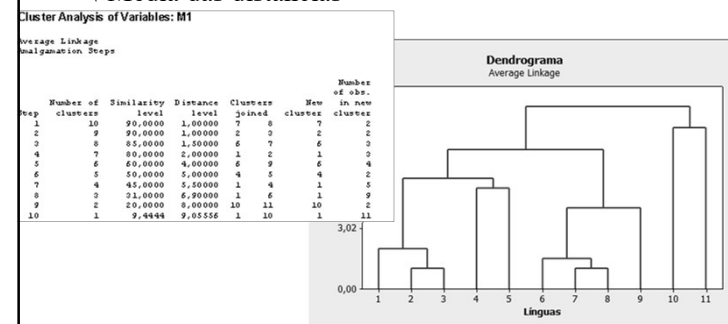
Análise Multivariada - 2022

71

**Exemplo 12.9**

• Numerais em 11 línguas

√ Média das distâncias



Análise Multivariada - 2022

73

- ✓ Leva a configuração muito semelhante ao dendrograma por ligação completa
- ✓ Combinações ocorrem em níveis diferentes
  - Distância é definida de maneira diferente

Análise Multivariada - 2022

74

### Exemplo 12.10

#### • Agrupamento de companhias (Média Distâncias)

✓ Dados de 22 concessionárias públicas (USA)

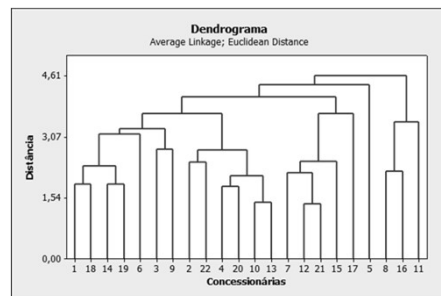
✓ Variáveis:

- $X_1$ : renda/dívidas
- $X_2$ : taxa de retorno de capitais
- $X_3$ : custo por capacidade instalada (kW)
- $X_4$ : fator de carga anual
- $X_5$ : pico de demanda (crescimento último ano)
- $X_6$ : Vendas (kWh por ano)
- $X_7$ : participação nucleares (%)
- $X_8$ : custo total de combustível (\$ por kWh)

✓ Dados: *BD\_multivariada.xls/public\_utilities*

Análise Multivariada - 2022

75



- ✓ Aplicada distância euclidiana a variáveis padronizadas
- ✓ Concessionárias tendem a se agrupar de acordo com localização geográfica

Análise Multivariada - 2022

76

#### • Método do Centróide:

✓ Distância entre dois grupos é definida como sendo a distância entre os vetores de médias (centróides)

- cada membro do par pertence a grupos diferentes



$$C_1 = \{X_1, X_2\} \text{ e } C_2 = \{X_3, X_4, X_5\}$$

$$d(C_1, C_2) = [(\bar{C}_1 - \bar{C}_2)'(\bar{C}_1 - \bar{C}_2)]^{1/2}$$

$$\bar{C}_1 = \frac{1}{2}(X_1 + X_2)$$

$$\bar{C}_2 = \frac{1}{3}(X_3 + X_4 + X_5)$$

Distância Euclidiana  
entre os dois grupos

Análise Multivariada - 2022

78



- É método direto e simples, mas em cada passo é necessário retornar aos dados originais para o cálculo da matriz de distâncias
  - √ exige mais tempo computacional
- Não pode ser usado em situações em que se dispõe apenas da matriz de distâncias entre os  $n$  elementos amostrais
  - √ Ao contrário dos métodos simples, complete e average linkage
- Quanto maior a quantidade de elementos amostrais ( $n$ ) e de variáveis ( $p$ ), menor a chance de empates entre valores da matriz de distâncias

Análise Multivariada - 2022

79

## Exemplo

- Dados 6 indivíduos de uma comunidade:
  - √ Renda (em salários mínimos)
  - √ Idade
  - √ Dados: (Fonte: Mingoti, 2005)

Indivíduo	Renda	Idade
A	9,60	28
B	8,40	31
C	2,40	42
D	18,20	38
E	3,90	25
F	6,40	41

- √ Agrupamento pelo método do centróide

Análise Multivariada - 2022

80

- Matriz de distâncias Euclidianas:

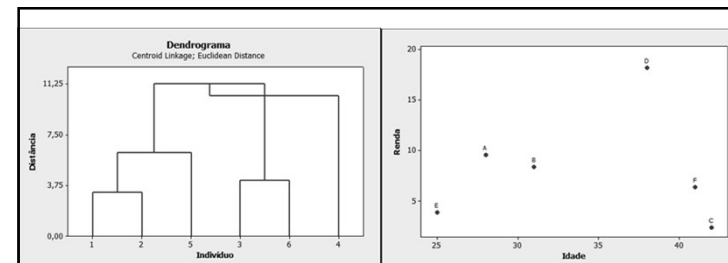
Matrix M3						
0,0000	3,2311	15,7429	13,1894	6,4413	13,3881	
3,2311	0,0000	12,5300	12,0433	7,5000	10,1980	
15,7429	12,5300	0,0000	16,2985	17,0660	4,1231	
13,1894	12,0433	16,2985	0,0000	19,3259	12,1754	
6,4413	7,5000	17,0660	19,3259	0,0000	16,1941	
13,3881	10,1980	4,1231	12,1754	16,1941	0,0000	

- Histórico do agrupamento:

Passo	g	Fusão	Distância (nível)
1	5	{A} e {B}	3,23
2	4	{C} e {F}	4,12
3	3	{A,B} e {E}	6,80
4	2	{A,B,E} e {C,F}	13,81
5	1	{A,B,E,C,F} e {D}	12,91

Análise Multivariada - 2022

81



- √ o nível de fusão do passo 5 foi menor que do passo 4
- √ Isso poderá ocorrer no método do centróide quando, em algum passo do algoritmo, houver empates entre valores da matriz de distâncias  $D$
- √ Quanto maior for o tamanho amostral e de variáveis, menor será a chance de ocorrência desta situação

Análise Multivariada - 2022

82

### Método de Ward

- Objetivo do procedimento:
  - ✓ Minimizar a perda de informação ao juntar 2 grupos
- Partição desejada:
  - ✓ A que produz os grupos mais heterogêneos entre si, com elementos homogêneos dentro de cada grupo
- Fundamento do método:
  - ✓ Em cada passo do agrupamento há mudança de variação entre os grupos e dentro dos grupos
  - ✓ Procedimento também denominado de mínima variância

Análise Multivariada - 2022

83

### Métodos anteriores:

- ✓ quando se passa de  $(n - k)$  para  $(n - k - 1)$  grupos o nível de fusão aumenta (nível de similaridade decresce) e a qualidade da partição decresce.
- ✓ Variação entre grupos diminui e a variação dentro dos grupos a

Análise Multivariada - 2022

84

### Procedimento

1. Cada elemento é considerado um único *cluster*;
2. Em cada passo calcula-se a soma da distância Euclidiana dentro dos grupos:

$$SSR = \sum_{i=1}^{g_k} SS_i$$

- ✓ SSR: soma dos quadrados total (dentro) dos grupos
- ✓  $g_k$ : número de grupos no passo  $k$
- ✓  $SS_i$ : soma dos quadrados do cluster  $i$

Análise Multivariada - 2022

85

$$SS_i = \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_i)'(X_{ij} - \bar{X}_i)$$

- ✓  $SS_i$ : soma dos quadrados do cluster  $i$
- ✓  $n_i$ : quantidade de elementos do *cluster*  $C_i$  (passo  $k$ )
- ✓  $X_{ij}$ : vetor de observações do  $j$ -ésimo elemento amostral que pertence ao  $i$ -ésimo conglomerado
- ✓  $\bar{X}_i$ : centróide do *cluster*  $i$

Análise Multivariada - 2022

86

3. Em cada passo do algoritmo, combinam-se os dois conglomerados que minimizam a distância entre os conglomerados  $C_i$  e  $C_j$ , definida como:

$$d(C_i, C_j) = \left[ \frac{n_i n_j}{n_i + n_j} \right] (\bar{X}_i - \bar{X}_j)' (\bar{X}_i - \bar{X}_j)$$

- ✓  $d(C_i, C_j)$  é a soma de quadrados entre os clusters  $C_i$  e  $C_j$

Análise Multivariada - 2022

87

• Comentários:

- ✓ Em cada passo, o método combina os dois conglomerados que resultam no menor valor de SSR
- ✓ Prova-se que  $d(C_i, C_j)$  é a diferença entre o valor de SSR depois e antes de se combinar os clusters  $C_i$  e  $C_j$ .
- ✓ Os métodos de Ward e do centróide usam o vetor de médias amostrais como representantes da informação global dos conglomerados em cada passo do processo de agrupamento
- ✓ A distância considera a diferença dos tamanhos dos conglomerados na comparação
  - $n_i n_j / (n_i + n_j)$  penalizam as comparações (maiores grupos → maiores distâncias)

Análise Multivariada - 2022

88

- O método do centróide não tem qualquer ponderação em relação ao tamanho dos *clusters*
- Para usar o método de Ward basta que as variáveis sejam quantitativas
  - ✓ Para o cálculo do vetor de médias
  - ✓ Não depende de se conhecer a distribuição da população
- Sob certas condições, há uma relação entre o método de Ward e o método de máxima verossimilhança
  - ✓ Se a distribuição das variáveis for normal p-variada

Análise Multivariada - 2022

89

- O método de Ward baseia-se na noção de que espera-se que os *clusters* de observações multivariadas tenham forma aproximadamente elíptica
- É um precursor de métodos de aglomeração não-hierárquicos
  - ✓ Otimizam algum critério para dividir os dados em um número determinado de grupos elípticos

Análise Multivariada - 2022

90

### Exemplo

- Dados 6 indivíduos de uma comunidade:
  - ✓ Renda (em salários mínimos)
  - ✓ Idade
  - ✓ Dados: (Fonte: Mingoti, 2005)

Indivíduo	Renda	Idade
A	9,60	28
B	8,40	31
C	2,40	42
D	18,20	38
E	3,90	25
F	6,40	41

✓ Agrupamento pelo método de Ward

Análise Multivariada - 2022

91

- Matriz de distâncias Euclidianas:

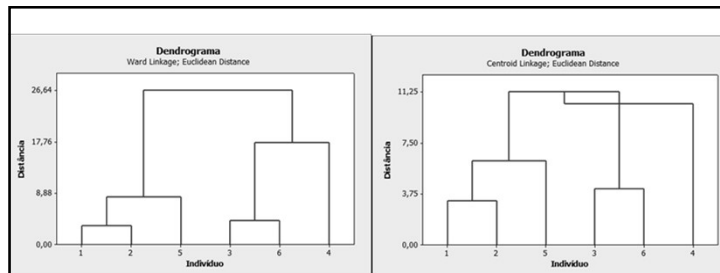
Matrix M3						
0,0000	3,2311	15,7429	13,1894	6,4413	13,3881	
3,2311	0,0000	12,5300	12,0433	7,5000	10,1980	
15,7429	12,5300	0,0000	16,2985	17,0660	4,1231	
13,1894	12,0433	16,2985	0,0000	19,3259	12,1754	
6,4413	7,5000	17,0660	19,3259	0,0000	16,1941	
13,3881	10,1980	4,1231	12,1754	16,1941	0,0000	

- Histórico do agrupamento:

Passo	g	Fusão	Distância (nível)
1	5	{A} e {B}	3,23
2	4	{C} e {F}	4,12
3	3	{A,B} e {E}	8,21
4	2	{C,F} e {D}	17,61
5	1	{A,B,E} e {C,F,D}	26,64

Análise Multivariada - 2022

92



- ✓ Os grupos finais foram os mesmos obtidos com o método do centróide
- ✓ Não houve inversão

Análise Multivariada - 2022

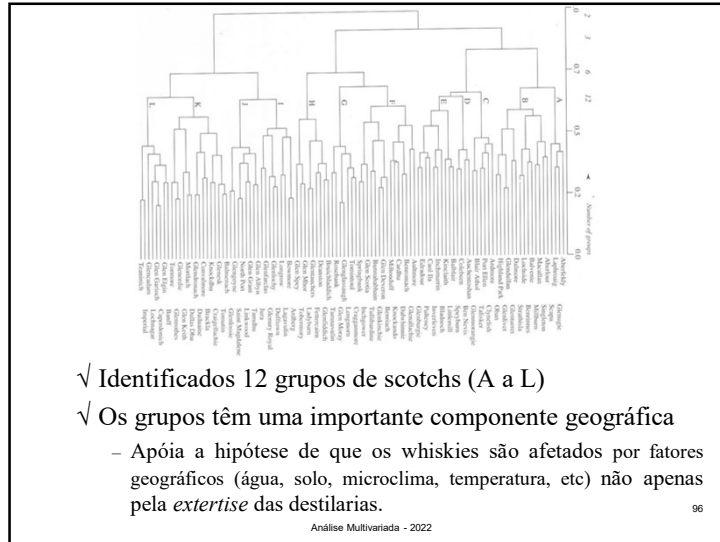
93

### Exemplo 12.11 – Pure Malt

- Agrupamento de 109 marcas de *scotch* de diferentes destilarias
- 68 variáveis binárias para medir as características do whiskey
  - ✓ relacionadas com: cor, corpo, aroma, etc.
- Objetivos:
  - ✓ Determinar os principais tipos de whiskies
  - ✓ Determinar suas principais características
  - ✓ Saber se os grupos correspondem a diferentes regiões
    - são afetados por solo, temperatura, condições da água
- Variáveis binárias são escaladas

Análise Multivariada - 2022

95



## Métodos Hierárquicos – Comentários Finais

- Fontes de erros e de variação não são formalmente considerados nos procedimentos hierárquicos
  - ✓ Significa que esses métodos são sensíveis a *outliers* ou pontos de perturbação
- Deve-se sempre verificar a sensibilidade da configuração dos grupos
  - ✓ Os métodos não permitem a realocação de objetos que possam ter sido agrupados incorretamente nos estágios iniciais

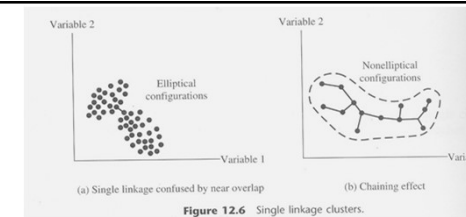
Análise Multivariada - 2022

97

- É recomendado tentar vários métodos de agrupamento e de atribuição de distâncias (similaridades)
- Empates na matriz de distâncias podem produzir múltiplas soluções ao problema de agrupamento hierárquico
- A maioria dos métodos produz clusters esféricos ou elípticos

Análise Multivariada - 2022

98



- O método de ligação simples é um dos poucos métodos que pode delinear cluster não-elípticos
  - ✓ Tem a capacidade de gerar estruturas geométricas diferentes
  - ✓ Tem a tendência de formar strings longas (chaining)
  - ✓ Entretanto, ele é incapaz de perceber grupos pouco separados

Análise Multivariada - 2022

99

- Os clusters formados pelo método de ligação simples não serão modificados por qualquer atribuição de distância (similaridade) que dá as mesmas ordenações relativas
  - ✓ Em particular, qualquer um dos coeficientes de similaridade monotônicos (Tabela 12.2)
- O método de ligação completa tende a produzir conglomerados de aproximadamente mesmo diâmetro
  - ✓ Tem a tendência de isolar os valores discrepantes nos estágios iniciais do agrupamento

Análise Multivariada - 2022

100

- O método da média das distâncias tende a produzir conglomerados de aproximadamente mesma variância interna
  - ✓ Em geral, produz melhores partições que os métodos de ligação simples e completa
- Os métodos de ligação simples, completa e da média podem ser utilizados tanto para variáveis quantitativas quanto para variáveis qualitativas
- Os métodos do centróide e de Ward são apropriados apenas para variáveis quantitativas

Análise Multivariada - 2022

101

- O método de Ward tende a produzir grupos com aproximadamente o mesmo número de elementos e tem como base principal os princípios de análise de variância
- Com um número maior de dados amostrais ( $n$ ) ou de variáveis ( $p$ ), necessariamente não irá ocorrer a igualdade das soluções apresentadas pelos vários métodos
  - ✓ Espera-se sempre que haja uma certa consistência entre as soluções obtidas por métodos diferentes

Análise Multivariada - 2022

102

### **Métodos para Encontrar o Número $g$ de Clusters da Partição Final**

- Problema de agrupamento:
  - ✓ Como escolher o número final ( $g$ ) de grupos que define a partição do conjunto de dados?
  - ✓ Qual o passo  $k$  em que o algoritmo de agrupamento deve ser interrompido?

Análise Multivariada - 2022

103

• Critério 1 – Análise do comportamento do nível de fusão (distância)

- ✓ À medida que o algoritmo avança, a similaridade dos grupos diminui (distância aumenta)
- ✓ Gráfico do passo (ou número de grupos) vs. nível de distância (nível de fusão)
  - Verifica-se a existência de “saltos” relativamente grandes
    - Pontos indicadores do momento ideal de parada (número final de conglomerados)
    - Se observados vários pontos de “saltos” pode-se delimitar uma região de prováveis valores do número de grupos  $g$  (deve ser investigado por outro procedimento)
- ✓ Pode-se usar o dendrograma quando  $n$  não for muito grande

Análise Multivariada - 2022

104

• Critério 2 – Análise do comportamento do nível de similaridade

- ✓ Similar ao critério 1
  - Observa-se o nível de similaridade (ao invés da distância)
- ✓ Nível de similaridade:
 
$$S_{ij} = \left( 1 - \frac{d_{ij}}{\max(d_{rs}), r, s = 1, \dots, n} \right) \times 100\%$$
  - $\max(d_{rs})$ : maior distância entre os  $n$  elementos amostrais na matriz de distâncias  $D_{n \times n}$  do início do processamento

Análise Multivariada - 2022

105

- ✓ Procura-se detectar pontos em que haja um decrescimento acentuado na similaridade dos conglomerados unidos
  - indicam a interrupção do algoritmo de agrupamento
  - número final de *clusters* ( $g$ ) está relacionado com o estágio em que o algoritmo foi interrompido
- ✓ Em geral, a escolha de valores de similaridade acima de 90% leva a um número de grupos muito elevado

Análise Multivariada - 2022

106

• Critério 3 – Análise da soma dos quadrados entre grupos:  $R^2$

- ✓ É possível calcular a soma de quadrados **entre clusters** e **dentro** dos grupos, em cada passo do procedimento
  - ✓ Em partição com  $g^*$  grupos, sejam:
    - $\mathbf{X}'_{ij} = (X_{i1,j}, X_{i2,j}, \dots, X_{ip,j})$
    - vetor de medidas observadas para o  $j$ -ésimo elemento amostral do  $i$ -ésimo grupo
    - $\bar{\mathbf{X}}'_i = (\bar{X}_{i1}, \bar{X}_{i2}, \dots, \bar{X}_{ip})$
    - vetor de médias do  $i$ -ésimo grupo (sem considerar partição)
    - $\bar{\mathbf{X}}' = (\bar{X}_{.1}, \bar{X}_{.2}, \dots, \bar{X}_{.p})$
- $$\bar{X}_{.r} = \frac{1}{n} \sum_{i=1}^{g^*} \sum_{j=1}^{n_i} X_{ijr}, \quad r = 1, 2, \dots, p$$

Análise Multivariada - 2022

107

- ✓ Soma dos quadrados total corrigida para a média global em cada variável

$$SST_C = \sum_{i=1}^{g^*} \sum_{j=1}^{n_i} (X_{ij} - \bar{X})(X_{ij} - \bar{X})$$

- ✓ Soma dos quadrados total dentro dos grupos da partição

$$SSR = \sum_{i=1}^{g^*} SS_i = \sum_{i=1}^{g^*} \sum_{j=1}^{n_i} (X_{ij} - \bar{X}_{i.})(X_{ij} - \bar{X}_{i.})$$

- ✓ Soma dos quadrados total entre os  $g^*$  grupos



Análise Multivariada - 2022

108

- ✓ Coeficiente  $R^2$  da partição:  $R^2 = \frac{SSB}{SST_C}$

- ✓ Quanto maior o valor de  $R^2$ , maior será a soma de quadrados entre grupos e menor será a soma de quadrados residual SSR

- ✓ Procedimento para escolha de  $g$

- Gráfico do passo do agrupamento vs.  $R^2$
- Procurar identificar algum ponto de ‘salto’ relativamente grande em relação aos demais
  - indica momento ideal da parada
- Gráfico é sempre decrescente
- maior valor de  $g^*$ , menor a variabilidade interna dos grupos e maior será o valor de  $R^2$ 
  - máximo  $R^2 = 1$  (para  $g^* = n$ )

Análise Multivariada - 2022

109

#### • Estratégia:

- ✓ Definir uma região de valores plausíveis para o número de grupos  $g$
- ✓ Utilizar o critério 3 dentro da região estabelecida

Análise Multivariada - 2022

110

#### • Critério 4 – Estatística Pseudo F

- ✓ (Calniski e Harabasz, 1974)

- ✓ Calcular estatística  $F$  em cada passo do agrupamento

$$F = \frac{\frac{SSB}{(g^*-1)}}{\frac{SSR}{(n-g^*)}} = \frac{(n-g^*)}{(g^*-1)} \frac{R^2}{1-R^2}$$

$g^*$ : número de grupos da partição em análise

- ✓ Idéia do teste:

- Em cada passo do agrupamento estaria sendo feito um teste F de análise de variância

- ✓ Importante:

- Na prática, não ocorre alocação aleatória
- A maioria dos métodos usa métodos de agrupamento com critérios determinísticos para partição dos dados

Análise Multivariada - 2022

111



- ✓ Se os elementos amostrais são provenientes de uma distribuição normal p-variada e quando os elementos são alocados aleatoriamente nos grupos
- ✓  $F \sim F_{p(g^* - 1), p(n - g^*)}$
- ✓ Se  $F$  é monotonicamente crescente com  $g^*$ , os dados sugerem que não há qualquer estrutura 'natural' de partição dos dados
- ✓ Se  $F$  apresentar um valor máximo, o número de conglomerados corresponderá à partição 'ideal'

Análise Multivariada - 2022

112

- ✓ Busca-se o maior valor de  $F$ 
  - Busca-se partição com maior heterogeneidade dos grupos
  - valor relacionado com a menor probabilidade de significância do teste
  - Estaria rejeitando a igualdade de vetores de médias populacionais com maior significância

Análise Multivariada - 2022

113

### • Critério 5 – Correlação Semiparcial (Método de Ward)

- ✓ Em determinado passo,  $C_k = C_i \cup C_j$

$$SPR^2 = \frac{B_{ij}}{SST_C}$$

Coefficiente de correlação parcial da partição

$$B_{ij} = \frac{n_i n_j}{n_i + n_j} (\bar{X}_{i\cdot} - \bar{X}_{j\cdot})' (\bar{X}_{i\cdot} - \bar{X}_{j\cdot})$$

Distância entre grupos – Método de Ward

1. Calcula-se  $SPR^2$  em cada passo
2. Gráfico passo vs.  $SPR^2$
3. Busca-se no gráfico salto consideravelmente maior que os restantes
4. Ponto indica partição ideal (parada do algoritmo de agrupamento)

Análise Multivariada - 2022

114

- ✓ A função  $SPR^2$  é não decrescente
- ✓ Se o agrupamento dos dados foi feito pelo método de Ward, o critério do coeficiente de correlação semiparcial equivalerá à aplicação do critério 1.

Análise Multivariada - 2022

115

• **Crítério 6 – Estatística Pseudo T<sup>2</sup>**

✓ Em determinado passo,  $C_k = C_i \cup C_j$

$$T^2 = \frac{B_{ij}}{\left[ \sum_{r \in C_i} \|\mathbf{X}_{ir} - \bar{\mathbf{X}}_i\|^2 + \sum_{r \in C_j} \|\mathbf{X}_{jr} - \bar{\mathbf{X}}_j\|^2 \right] (n_i + n_j)^{-1}}$$

$$\|\mathbf{X}_{sr} - \bar{\mathbf{X}}_{s\cdot}\| = (\mathbf{X}_{sr} - \bar{\mathbf{X}}_{s\cdot})'(\mathbf{X}_{sr} - \bar{\mathbf{X}}_{s\cdot}), s = i, j$$

✓ Sob as suposições de normalidade p-variada e alocação aleatória dos grupos

$$T^2 \sim F_{p, (n_i+n_j-2)}$$

✓ Na prática, não se tem alocação aleatória dos grupos

✓ Ideia do teste:

- Teste de comparação de média de dois grupos, unidos para formar novo grupo

116

Análise Multivariada - 2022

✓ Gráfico passo vs. valor da Pseudo T<sup>2</sup>

- Busca-se no gráfico o valor máximo

✓ O valor de g correspondente ao máximo (ou aquele imediatamente anterior) é escolhido como o número provável de grupos da partição final

✓ Busca-se o maior valor de T<sup>2</sup>

- aquele relacionado com a menor probabilidade de significância

(Rejeita a igualdade dos vetores de média com maior significância)

- Se a igualdade entre os vetores de médias é rejeitada, os dois clusters deveriam ser unidos para formar um único agrupamento

117

Análise Multivariada - 2022

• **Crítério 6 – Estatística CCC (*Cubic Clustering Criterium*)**

✓ Sarle (1983)

✓ Compara-se o valor esperado do coeficiente R<sup>2</sup> com a aproximação do valor esperado de r<sup>2</sup> sob a suposição de que os grupos são gerados de acordo com uma uniforme p-dimensional

✓ CCC indicaria a presença de estrutura de agrupamento diferente da partição uniforme

✓ A quantidade de grupos da partição final estaria relacionada com valores de CCC > 3

✓ Está implementada no software estatístico SAS

118

Análise Multivariada - 2022

**Exemplo 6.8**

Mingoti, 2005

- Dados relativos a 21 países (ONU, 2002)
- Variáveis:
  - ✓ Expectativa de vida
  - ✓ Educação
  - ✓ Renda (PIB)
  - ✓ Estabilidade política e de segurança
- Método de agrupamento: Ward
- Conjunto de dados: *BD\_multivariada.xls/paises*

120

Análise Multivariada - 2022

• Minitab

Cluster Analysis of Observations: Índice de Ex; Índice de Ed; Índice PIB; ...

Squared Euclidean Distance, Ward Linkage  
Amalgamation Steps

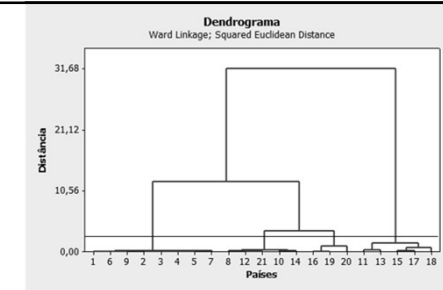
Step	Number of clusters	Similarity level	Distance level	Clusters joined	New cluster	Number of obs. in new cluster
1	20	99,995	0,0006	2 3	2	2
2	19	99,966	0,0042	1 6	1	2
3	18	99,904	0,0091	4 5	4	2
4	17	99,923	0,0095	2 4	2	4
5	16	99,821	0,0219	12 21	12	2
6	15	99,613	0,0470	1 9	1	3
7	14	99,510	0,0598	16 19	16	2
8	13	99,462	0,0657	8 12	8	3
9	12	99,377	0,0761	2 7	2	5
10	11	98,999	0,1221	10 14	10	2
11	10	98,963	0,1266	15 17	15	2
12	9	98,622	0,1691	1 2	1	8
13	8	97,997	0,2445	11 13	11	2
14	7	97,535	0,3008	8 10	8	5
15	6	94,602	0,6597	15 18	15	3
16	5	92,489	0,9166	16 20	16	3
17	4	88,114	1,4505	11 15	11	5
18	3	71,202	3,5145	8 16	8	8
19	2	1,220	12,055	1 8	1	16
20	1	-159,594	31,6803	1 11	1	21

Final Partition  
Number of clusters: 1

	Number of observations	cluster sum of squares	Average distance from centroid	Maximum distance from centroid
Cluster1	21	25,7654	0,987665	2,3069

Análise Multivariada - 2022

121



- ✓ Visualmente, é razoável definir 4 grupos de países
- ✓ Analisar queda de similaridade entre os passos 16 e 18

Análise Multivariada - 2022

122

• Resultados da análise de agrupamento:

Passo	g*	Similaridade	Distância	R <sup>2</sup>	Pseudo F	sp R <sup>2</sup>	Pseudo T <sup>2</sup>	CCC
1	20	99,99	0,001	1,000	4520,0	0,0000		
2	19	99,97	0,004	1,000	1193,0	0,0001		
3	18	99,93	0,008	1,000	705,0	0,0002		
4	17	99,92	0,009	1,000	576,0	0,0002	2,2	
5	16	99,82	0,022	0,999	388,0	0,0004		
6	15	99,61	0,047	0,998	241,0	0,0009	11,3	
7	14	99,51	0,060	0,997	183,0	0,0012		
8	13	99,46	0,066	0,996	158,0	0,0013	3,0	
9	12	99,38	0,076	0,994	143,0	0,0015	12,6	
10	11	99,00	0,122	0,992	123,0	0,0024		
11	10	98,96	0,127	0,989	115,0	0,0025		
12	9	98,62	0,168	0,986	107,0	0,0033	6,9	
13	8	98,00	0,245	0,981	96,4	0,0047		
14	7	97,54	0,301	0,976	93,5	0,0058	4,3	
15	6	94,60	0,659	0,963	77,8	0,0126	5,2	
16	5	92,49	0,917	0,945	68,8	0,0178	15,3	
17	4	88,11	1,450	0,917	62,5	0,0281	4,2	-0,12
18	3	71,20	3,514	0,849	50,5	0,0682	14,2	-0,65
19	2	1,22	12,055	0,815	30,3	0,2339	31,8	-1,80
20	1	-159,59	31,680	0,000		0,6148	30,3	

✓ Do passo 17 para 18:

- Perda mais acentuada de similaridade
- O valor de R<sup>2</sup> passa de 0,917 para 0,849
- Valores da Pseudo F e do CCC decrescem substancialmente
- Pseudo T<sup>2</sup> e SPR<sup>2</sup> crescem acentuadamente

Análise Multivariada - 2022

123

• Medidas descritivas dos grupos formados:

Grupos (SQ)	Países	Média			
		Expectativa de vida	Educação	PIB	Estabilidade política
1 (0,157) n <sub>1</sub> = 8	Austrália, Canadá, Cingapura, Estados Unidos, França, Japão, Reino Unido, Uruguai	0,8838	0,9538	0,9075	1,1850
2 (0,255) n <sub>2</sub> = 5	Argentina, Brasil, China, Cuba, Egito	0,7660	0,8140	0,6740	0,3380
3 (1,240) n <sub>3</sub> = 5	Angola, Colômbia, Nigéria, Paraguai, Serra Leoa	0,5060	0,5900	0,4940	-1,3660
4 (0,488) n <sub>4</sub> = 3	Etiópia, Moçambique, Senegal	0,3400	0,3633	0,3767	-0,3433
Global n = 21	Todos	0,6881	0,7495	0,6776	0,1580

✓ Grupo 1 – ‘Primeiro Mundo’

- países com maiores índice de desenvolvimento

✓ Grupo 4 – alguns países africanos

- menores índices em todas as variáveis

Análise Multivariada - 2022

124

✓ Variável estabilidade política e segurança:

- Grande diferença de comportamento dos grupos 1 e 2 em relação aos grupos 3 e 4
- Grupo 1 é o de maior estabilidade e o grupo 3 de menor

✓ Dispersão interna é menor no grupo 1 e maior no grupo 3

Análise Multivariada - 2022

125

## Similaridade e Associação

- Em algumas aplicações, deseja-se agrupar as variáveis e não os itens

✓ Em geral, medidas de associação amostral são usadas como medida de similaridade entre variáveis

- Em algumas aplicações de agrupamento, substitui-se correlação negativa por seu valor absoluto

Análise Multivariada - 2022

127

## Técnicas Hierárquicas e Seleção de Variáveis

- Os métodos hierárquicos podem ser úteis na seleção das variáveis mais importantes na caracterização de determinada situação
- Métodos de ligação simples, completa e da média
  - ✓ É necessária apenas matriz inicial que represente proximidade (ou similaridade) entre os elementos amostrais
  - ✓ É necessário escolher uma matriz inicial que represente o relacionamento dessas variáveis
    - Interesse: agrupar as variáveis mais similares entre si (separar aquelas com informações diferenciadas)

Análise Multivariada - 2022

128

## Variáveis quantitativas:

- Pode-se usar coeficiente de correlação de Pearson
  - ✓ Expressa similaridade com relação à associação linear
  - ✓ Quanto maior seu valor absoluto, maior a aproximação entre as variáveis
- Matriz de correlação amostral não é uma matriz de distâncias (ou proximidades)
  - ✓ Transformação mais simples
$$D_{p \times p} = 1 - \text{Abs}(R_{p \times p})$$
- Podem ser usados coeficientes de correlação não paramétricos
  - ✓ Spearman, Kendall, etc.

Análise Multivariada - 2022

129

### Exemplo

- Matriz de correlação amostral (R):

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
X <sub>1</sub>	1					
X <sub>2</sub>	0,57	1				
X <sub>3</sub>	0,51	0,60	1			
X <sub>4</sub>	0,39	0,38	0,43	1		
X <sub>5</sub>	0,46	0,32	0,40	0,50	1	
X <sub>6</sub>	0,35	0,72	0,45	0,58	0,30	1

✓ X<sub>2</sub> e X<sub>6</sub> são mais similares ( $r_{26} = 0,72$ )

- $D_{6 \times 6} = 1 - \text{Abs}(R_{6 \times 6})$

	X <sub>1</sub>	X <sub>2</sub>	X <sub>3</sub>	X <sub>4</sub>	X <sub>5</sub>	X <sub>6</sub>
X <sub>1</sub>	0					
X <sub>2</sub>	0,43	0				
X <sub>3</sub>	0,49	0,40	0			
X <sub>4</sub>	0,61	0,62	0,57	0		
X <sub>5</sub>	0,54	0,68	0,60	0,50	0	
X <sub>6</sub>	0,65	0,28	0,55	0,42	0,70	0

Análise Multivariada - 2022

130

- Método de Ligação Simples

Passo	g	Fusão	Nível Fusão
1	5	X <sub>2</sub> e X <sub>6</sub>	0,28
2	4	X <sub>2</sub> , X <sub>6</sub> e X <sub>3</sub>	0,40
3	3	X <sub>2</sub> , X <sub>6</sub> , X <sub>3</sub> e X <sub>4</sub>	0,42
4	2	X <sub>2</sub> , X <sub>6</sub> , X <sub>3</sub> , X <sub>4</sub> e X <sub>1</sub>	0,43
5	1	X <sub>2</sub> , X <sub>6</sub> , X <sub>3</sub> , X <sub>4</sub> , X <sub>1</sub> e X <sub>5</sub>	0,50

✓ No passo 3

- $C_1 = \{X_2, X_6, X_3, X_4\}$
- $C_2 = \{X_1\}$
- $C_3 = \{X_5\}$

Suponha escolher 3 dentre as 6 variáveis:  
X<sub>1</sub>  
X<sub>5</sub>  
Uma das variáveis de C<sub>1</sub>

Análise Multivariada - 2022

131

- Medidas de similaridade para variáveis categóricas:

- ✓ Coeficiente qui-quadrado
- ✓ Coeficiente de contingência de Pearson
- ✓ Coeficiente de concordância de Kappa

- Outros Coeficientes:

- ✓ Podem-se desenvolver medidas de associação (similaridade) análogos aos coeficientes estabelecidos anteriormente (Tabela 12.2) – Troca-se  $p$  por  $n$ .

Análise Multivariada - 2022

132

- Variáveis Binárias:

- ✓ Os dados podem ser agrupados na forma de tabela de contingência
- ✓ Para cada par de variáveis, há  $n$  itens categorizados na tabela

		Variável $k$		
		1	0	Total
Variável $i$	1	a	b	a + b
	0	c	d	c + d
Total		a + c	b + d	n = a + b + c + d

Análise Multivariada - 2022

133

- Coeficiente de Correlação Phi

$$r = \frac{ad - bc}{[(a+b)(c+d)(a+c)(b+d)]^{1/2}}$$

✓ Pode ser tomado como medida de similaridade entre as duas variáveis

✓  $r$  está relacionado com a estatística  $\chi^2$  para teste de independência entre duas variáveis categóricas

$$r^2 = \frac{\chi^2}{n}$$

✓ Para  $n$  fixo, uma correlação (similaridade) grande é consistente com a ausência de independência

Análise Multivariada - 2022

135

## Exemplo 12.7

- Agrupamento de variáveis (Ligação Completa)

✓ Dados de 22 concessionárias públicas (USA)

✓ Variáveis:

- $X_1$ : renda/dívidas
- $X_2$ : taxa de retorno de capitais
- $X_3$ : custo por capacidade instalada (kW)
- $X_4$ : fator de carga anual
- $X_5$ : pico de demanda (crescimento último ano)
- $X_6$ : Vendas (kWh por ano)
- $X_7$ : participação nucleares (%)
- $X_8$ : custo total de combustível (\$ por kWh)

✓ Dados: *BD\_multivariada.xls/public\_utilities*

Análise Multivariada - 2022

136

- Coeficiente de correlação para medir similaridade

✓ variáveis com grandes correlações negativas são consideradas muito dissimilares

✓ variáveis com grandes correlações positivas são consideradas muito similares

✓ distância entre *clusters* é medida como menor similaridade entre grupos

- Matriz de correlações:

Correlations: X1; X2; X3; X4; X5; X6; X7; X8

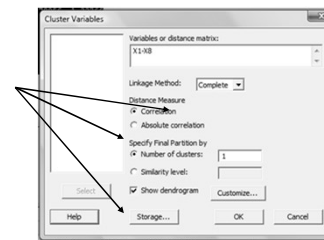
	X1	X2	X3	X4	X5	X6	X7
X2	0,643						
X3	-0,103	-0,348					
X4	-0,082	-0,086	0,100				
X5	-0,259	-0,260	0,435	0,033			
X6	-0,152	-0,010	0,028	-0,288	0,176		
X7	0,045	0,211	0,115	-0,164	-0,019	-0,374	
X8	-0,013	-0,328	0,005	0,486	-0,007	-0,561	-0,185

Análise Multivariada - 2022

137

- Minitab

Stat > Multivariate > Cluster Variables →



- Matriz de distâncias:  $D_{8 \times 8} = 1 - R_{8 \times 8}$

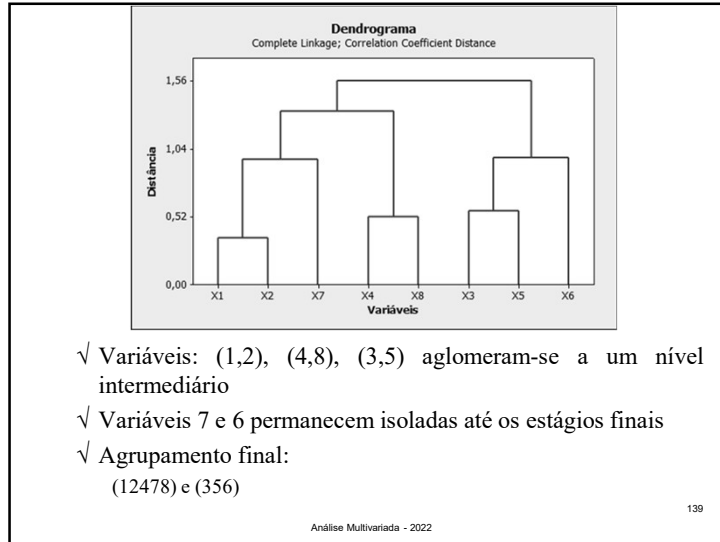
Data Display

Matriz de Distâncias

0,00000	0,35726	1,10279	1,08203	1,25911	1,15167	0,95520	1,01337
0,35726	0,00000	1,34755	1,08634	1,26011	1,00962	0,78856	1,32766
1,10279	1,34755	0,00000	0,89969	0,56463	0,97201	0,88534	0,99478
1,08203	1,08634	0,89969	0,00000	0,96652	1,28794	1,16416	0,51450
1,25911	1,26011	0,56463	0,96652	0,00000	0,62358	1,01913	1,00713
1,15167	1,00962	0,97201	1,28794	0,62358	0,00000	1,37369	1,56053
0,95520	0,78856	0,88534	1,16416	1,01913	1,37369	0,00000	1,18509
1,01337	1,32766	0,99478	0,51450	1,00713	1,56053	1,18509	0,00000

Análise Multivariada - 2022

138



## Comentários

- Há várias maneiras de medir similaridade entre pares de objetos:
  - ✓ distâncias (12-1 a 12-5)
  - ✓ Coeficientes (Tabela 12-2) – para agrupar itens
  - ✓ Correlações – para agrupar variáveis
- Podem ser usadas frequências

## Técnicas de Agrupamento Não Hierárquicas

- Objetivo:
  - ✓ Encontrar diretamente uma partição de  $n$  elementos em  $k$  grupos
  - ✓ Requisitos:
    - coesão interna (semelhança interna)
    - isolamento (separação) dos clusters formados
- Busca da “melhor” partição de ordem  $k$ 
  - ✓ Satisfaz algum critério de qualidade
  - ✓ Procedimentos computacionais para investigar partições ‘quase’ ótima (inviável a busca exaustiva)

- Métodos Não Hierárquicos vs. Hierárquicos :
  - ✓ Especificação prévia do número de cluster (ao contrário das técnicas aglomerativas)
  - ✓ Novos grupos podem ser formados pela divisão (ou junção) de grupos já combinados:
    - Se em um passo do algoritmo, dois elementos tiverem sido colocados em um mesmo grupo, não significa que estarão juntos na partição final
    - Não é mais possível a construção de dendrogramas
  - ✓ Em geral, são do tipo iterativo

- ✓ Tem maior capacidade de analisar grande número de dados
- ✓ A matriz de distância não tem de ser calculada e os dados básicos não precisam ser armazenados durante a execução do procedimento
- ✓ Métodos hierárquicos são mais adequados para agrupar itens que variáveis

Análise Multivariada - 2022

143

### Métodos Não Hierárquicos – Estrutura

- Iniciam-se:
  1. partição inicial de itens em grupos
  2. conjunto inicial de sementes que formarão o núcleo dos clusters
- Escolha das configurações iniciais pode afetar partição final
  - ✓ Viés na escolha das sementes iniciais
  - ✓ Alternativas:
    - Seleção aleatória de sementes
    - Partição aleatória de itens em grupos iniciais

Análise Multivariada - 2022

144

### Métodos Não Hierárquicos – Procedimentos

Alguns procedimentos:

- Método das  $k$ -Médias (*k-Means*)
  - ✓ mais conhecido e popular
- Método Fuzzy  $c$ -Médias
- Redes Neurais Artificiais

Análise Multivariada - 2022

145

### Método das $k$ -Médias

- Provavelmente, um dos mais conhecidos e mais utilizados
- Idéia Básica:
  - ✓ Cada elemento amostral é alocado àquele *cluster* cujo centróide é o mais próximo do elemento

Análise Multivariada - 2022

146



### Passos do Procedimento

1. Escolhem-se  $k$  centróides para inicializar o processo de partição
  - ✓ Sementes ou protótipos
2. Cada elemento do conjunto de dados é comparado com cada centróide inicial
  - ✓ Alocação ao centróide menos distante
  - ✓ Em geral, utiliza-se distância Euclidiana
  - ✓ Aplicação a todos os  $n$  elementos amostrais

Análise Multivariada - 2022

147

3. Cálculo dos novos centróides para cada grupo formado no passo (2)
  - ✓ Repetição do passo (2), considerando os novos valores dos centróides
4. Os passos (2) e (3) são repetidos até que todos os elementos amostrais estejam “bem alocados” em seus grupos
  - ✓ “Bem alocados” = não é necessária realocação de elementos

Análise Multivariada - 2022

148

### Exemplo 12.12

- Agrupamento pelo Método das  $k$ -Médias:
  - ✓ Medidas das variáveis  $X_1$  e  $X_2$ :

Observações		
Item	$x_1$	$x_2$
A	5	3
B	-1	1
C	1	-2
D	-3	-2

- ✓ Dividir em  $k = 2$  grupos de maneira que os itens de um *cluster* sejam os mais próximos um dos outros e que estejam distantes em *clusters* diferentes

Análise Multivariada - 2022

149

### Implementação:

- ✓ Partição arbitrária em 2 *clusters*: (AB) e (CD)
- ✓ Cálculo das coordenadas  $(\bar{x}_1, \bar{x}_2)$  dos centróides:

Cluster	$\bar{x}_1$	$\bar{x}_2$
AB	2	2
CD	-1	-2

- ✓ Distância euclidiana de cada item

	AB	CD
A	$d(A,AB) = (5-2)^2 + (3-2)^2 = 10$	$d(A,CD) = (5+1)^2 + (3+2)^2 = 61$
B	$d(B,AB) = (-1-2)^2 + (1-2)^2 = 10$	$d(B,CD) = (-1+1)^2 + (1+2)^2 = 9$
C	$d(C,AB) = (1-2)^2 + (-2-2)^2 = 17$	$d(C,CD) = (1+1)^2 + (-2+2)^2 = 4$
D	$d(D,AB) = (-3-2)^2 + (-2-2)^2 = 41$	$d(D,CD) = (-3+1)^2 + (-2+2)^2 = 4$

B é agrupado ao *cluster* (CD)

Análise Multivariada - 2022

150

✓ Cálculo das coordenadas  $(\bar{x}_1, \bar{x}_2)$  dos centróides:

Cluster	$\bar{x}_1$	$\bar{x}_2$
A	5	3
BCD	-1	-1

✓ Distância euclidiana de cada item

	A	BCD
A	$d(A,A) = (5-5)^2 + (3-3)^2 = 0$	$d(A,BCD) = (5+1)^2 + (3+1)^2 = 52$
B	$d(B,A) = (-1-5)^2 + (1-3)^2 = 40$	$d(B,BCD) = (-1+1)^2 + (1+1)^2 = 4$
C	$d(C,A) = (1-5)^2 + (-2-3)^2 = 41$	$d(C,BCD) = (1+1)^2 + (-2+1)^2 = 5$
D	$d(D,A) = (-3-5)^2 + (-2-3)^2 = 89$	$d(D,BCD) = (-3+1)^2 + (-2+1)^2 = 5$

✓ O agrupamento se mantém e o processo pára

Análise Multivariada - 2022

151

• Agrupamento Final ( $k = 2$ )

✓ A e (BCD)

• Comentários:

✓ Para verificar a estabilidade da solução é recomendável reiniciar o algoritmo com uma nova partição inicial

✓ Uma tabela de centróides e das variâncias dentro dos grupos auxilia a delinear as diferenças entre os grupos

Análise Multivariada - 2022

152

## Sugestões para Escolha Cuidadosa das Sementes

• Sugestão 1: Uso de técnicas hierárquicas aglomerativas:

✓ Utiliza-se algum método de agrupamento hierárquico para se obter os  $k$  grupos iniciais

✓ Calcula-se o vetor de médias de cada grupo

✓ Esses vetores são usados como sementes iniciais

Análise Multivariada - 2022

154

• Sugestão 2: Escolha aleatória:

✓ As  $k$  sementes iniciais são escolhidas aleatoriamente dentro do conjunto de dados

✓ Sugestão amostragem aleatória simples sem reposição  
(estratégia simples, mas sem eficiência)

✓ Melhoria de eficiência na escolha:

- Selecionar  $m$  amostras aleatórias, constituídas de  $k$  sementes
- Cálculo do vetor de médias das  $k$  sementes selecionadas para cada grupo
- Esses vetores constituem os centróides de inicialização do processo de agrupamento das  $k$ -médias

Análise Multivariada - 2022

155

- Sugestão 3: Escolha por meio de uma variável aleatória:

- ✓ Escolhe-se uma variável aleatória dentre as  $p$  componentes em consideração
  - a variável por si só já induz uma certa “partição natural” dos dados
- ✓ Divide-se o domínio da variável em  $k$  intervalos
- ✓ A semente inicial será o centróide de cada intervalo

Análise Multivariada - 2022

156

- Sugestão 4: Observação dos valores discrepantes do conjunto de dados

- ✓ Análise estatística para buscar  $k$  elementos discrepantes no conjunto de dados
  - Discrepância em relação às  $p$  variáveis observadas
- ✓ Cada um desses elementos será a semente

Análise Multivariada - 2022

157

- Sugestão 5: Escolha prefixada

- ✓ Método não muito recomendável, pois, tem um alto grau de subjetividade
- ✓ Sementes escolhidas arbitrariamente
- ✓ Pode ser usadas em casos em há grande conhecimento do problema
  - buca-se validar solução já existente

Análise Multivariada - 2022

158

- Sugestão 6: Os  $k$  primeiros valores do banco de dados

- ✓ Usado como *default* pela maioria dos softwares
- ✓ Pode trazer bons resultados quando os  $k$  primeiros elementos amostrais são discrepantes entre si
  - (Não é recomendável quando são semelhantes)

Análise Multivariada - 2022

159

### Exemplo 7.1

Mingoti, 2005 – Continuação Ex. 6.8

- Dados relativos a 21 países (ONU, 2002)
- Variáveis:
  - ✓ Expectativa de vida
  - ✓ Educação
  - ✓ Renda (PIB)
  - ✓ Estabilidade política e de segurança
- Método de agrupamento:  $k$ -Médias
- Conjunto de dados: *BD\_multivariada.xls/paises*

Análise Multivariada - 2022

160

- Utiliza-se da Análise pelo Método de Ward:
  - ✓  $k = g = 4$  grupos para partição dos países
  - ✓ Sementes iniciais = centróides *clusters* finais
- Partição final:
  - ✓ a mesma obtida anteriormente

Grupos (SQ)	Países	Média			
		Expectativa de vida	Educação	PIB	Estabilidade política
1 (0,157) $n_1 = 8$	Austrália, Canadá, Cingapura, Estados Unidos, França, Japão Reino Unido, Uruguai	0,8838	0,9538	0,9075	1,1850
2 (0,255) $n_2 = 5$	Argentina, Brasil, China, Cuba, Egito	0,7660	0,8140	0,6740	0,3380
3 (1,240) $n_3 = 5$	Angola, Colômbia, Nigéria, Paraguai, Serra Leoa	0,5060	0,5900	0,4940	-1,3660
4 (0,488) $n_4 = 3$	Etiópia, Moçambique, Senegal	0,3400	0,3633	0,3767	-0,3433
Global $n = 21$	Todos	0,6881	0,7495	0,6776	0,1580

Análise Multivariada - 2022

161

- Sementes iniciais: Reino Unido, Brasil, Serra Leoa e Moçambique
  - ✓ Obtém-se mesma partição final

Análise Multivariada - 2022

162

- Sementes iniciais: 4 primeiros países do banco

Grupos (SQ)	Países	Média			
		Expectativa de vida	Educação	PIB	Estabilidade política
1 (0,091) $n_1 = 7$	Austrália, Canadá, Estados Unidos, França, Japão, Reino Unido, Uruguai	0,8843	0,9657	0,9071	1,1529
2 (0,748) $n_2 = 6$	Argentina, Brasil, China, Cuba, Egito, Moçambique	0,6783	0,7400	0,6271	0,3150
3 (2,188) $n_3 = 7$	Angola, Colômbia, Nigéria, Paraguai, Serra Leoa, Etiópia, Senegal	0,4729	0,5243	0,4626	-1,1514
4 (0,488) $n_4 = 1$	Cingapura	0,8800	0,8700	0,9100	1,4100
Global $n = 21$	Todos	0,6881	0,7495	0,6776	0,1580

Análise Multivariada - 2022

163

- ✓ Cingapura foi separada do *cluster* 1
- ✓ Moçambique deslocado para grupo do Brasil
- ✓ Grupo da Colômbia acrescido de Etiópia e Senegal

- ✓ Soma de quadrados dentro dos grupos:
  - Nova solução aumentou variabilidade dentro dos grupos 2 e 3

Análise Multivariada - 2022

164

## • Minitab:

**K-means Cluster Analysis: Índice de Ex; Índice de Ed; Índice PIB; Estabilidade**

Final Partition

Number of clusters: 4

	Number of observations	Within cluster sum of squares	Average distance from centroid	Maximum distance from centroid
Cluster1	7	2,188	0,528	0,856
Cluster2	1	0,000	0,000	0,000
Cluster3	7	0,091	0,105	0,204
Cluster4	6	0,748	0,308	0,641

Cluster Centroids

Variable	Cluster1	Cluster2	Cluster3	Cluster4
Índice de Expectativa de vida	0,4729	0,8800	0,8843	0,6783
Índice de Educação	0,5243	0,8700	0,9657	0,7400
Índice PIB	0,4629	0,9100	0,9071	0,6217
Estabilidade_politica_violência	-1,1514	1,4100	1,1529	0,3150

Variable	Grand centroid
Índice de Expectativa de vida	0,6881
Índice de Educação	0,7495
Índice PIB	0,6776
Estabilidade_politica_violência	0,1576

Distances Between Cluster Centroids

	Cluster1	Cluster2	Cluster3	Cluster4
Cluster1	0,0000	2,6545	2,4231	1,5048
Cluster2	2,6545	0,0000	0,2744	1,1575
Cluster3	2,4231	0,2744	0,0000	0,9364
Cluster4	1,5048	1,1575	0,9364	0,0000

165

- Soma dos Quadrados:
  - ✓ ANOVA para cada uma das  $p$  variáveis

Soma Quadrados Total Corrigida

$SST_c = 25,7654$

**One-way ANOVA: Índice de Expectativa de vida versus C35**

Source	DF	SS	MS	F	P
C35	3	0,6311	0,2104	6,83	0,003
Error	17	0,5234	0,0308		
Total	20	1,1545			

**One-way ANOVA: Índice de Educação versus C35**

Source	DF	SS	MS	F	P
C35	3	0,6974	0,2325	7,22	0,002
Error	17	0,5475	0,0322		
Total	20	1,2449			

**One-way ANOVA: Índice PIB versus C35**

Source	DF	SS	MS	F	P
C35	3	0,7644	0,2548	14,94	0,000
Error	17	0,2900	0,0171		
Total	20	1,0544			

**One-way ANOVA: Estabilidade\_politica\_violência versus C35**

Source	DF	SS	MS	F	P
C35	3	20,6458	6,8819	70,23	0,000
Error	17	1,6658	0,0980		
Total	20	22,3116			

Soma dos Quadrados Total dentro dos grupos  
 $SSR = 3,0267$

Análise Multivariada - 2022

166

## • Cálculo $R^2$ :

$$SSB = SST_C - SSR = 25,7654 - 3,0267 = 22,7387$$

$$R^2 = \frac{SSB}{SST_C} = \frac{22,7387}{25,7654} = 0,883$$

## • Comparação da qualidade das partições

Partição	Soma Quadrados Residual	Varição Residual Média	$R^2$
k-Médias (Ward)	2,1406	0,5351	0,917
k-Médias (*)	3,0267	0,7557	0,883

(\*) Os 4 primeiros elementos do banco de dados

✓ Utilização das sementes de Ward e 4 primeiros elementos do banco de dados como sementes

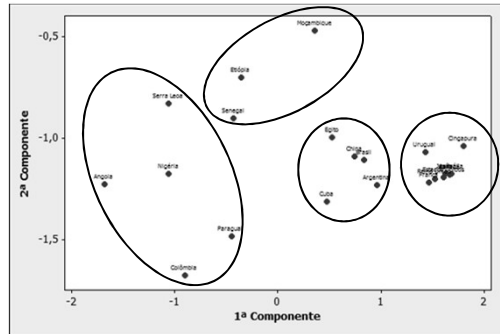
- Sementes de Ward levam a uma melhor solução de agrupamento pelo método das k-Médias.

Análise Multivariada - 2022

167

- Visualização espacial dos grupos:

✓ 2 primeiras componentes principais com base na matriz de covariâncias amostral



✓ É possível visualizar claramente os 4 grupos

–  $k$ -médias com sementes de Ward

168

Análise Multivariada - 2022

## Comentários Finais

- A escolha das sementes iniciais de agrupamento podem influenciar o agrupamento final

✓ Se duas ou mais sementes situarem-se em um único *cluster*, os grupos resultantes serão pouco diferenciados

✓ A existência de *outlier* pode produzir pelo menos um grupo com muitos itens dispersos

Análise Multivariada - 2022

169

- Há fortes argumentos para não se fixar o número de *clusters*  $k$

✓ Mesmo sabendo-se que a população consiste de  $k$  grupos, dependendo do método de amostragem, pode não aparecer na amostra os dados provenientes de um grupo mais raro

– Forçar  $k$  grupos levaria a *clusters* sem sentido

✓ Em casos em que o algoritmo requer o uso de um valor especificado de  $k$ , é sempre uma boa idéia executar novamente o algoritmo para diversas escolhas de  $k$

170

Análise Multivariada - 2022

## Escalonamento Multidimensional

### Escalonamento Multidimensional

- Método para representar dados multivariados em espaços de menor dimensão
- Objetivo básico:
  - ✓ Ajustar os dados originais em sistema coordenado de dimensões reduzidas
  - ✓ Dimensões são obtidas procurando preservar similaridades e dissimilaridades

Análise Multivariada - 2022

172

### Tipos de Escalonamento Multidimensional

- Métrico:
  - ✓ Situações nas quais p variáveis quantitativas são medidas em cada elemento amostral
- Não métrico:
  - ✓ Matriz de similaridade não é construída a partir de distâncias matemáticas
    - Julgamentos ou percepções dos indivíduos sobre os objetos
    - Em geral, dados provenientes de questionários sobre estímulos (produtos, marcas, serviços, etc.)

Análise Multivariada - 2022

173

### Método MDS – Escalonamento Métrico

- Seja uma matriz de distâncias (ou dissimilaridades)  $\mathbf{D}_{n \times n}$ .
  - ✓ Obtém-se uma matriz  $\mathbf{A}_{n \times n}$   $a_{ij} = -\frac{1}{2} [d_{ij}^2 - d_{i.}^2 - d_{.j}^2 + d_{..}^2]$ .
  - ✓ com  $d_{i.}^2 = \frac{1}{n} \sum_{j=1}^n d_{ij}^2$ ,  $d_{.j}^2 = \frac{1}{n} \sum_{i=1}^n d_{ij}^2$  e  $d_{..}^2 = \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^n d_{ij}^2$ .
  - ✓ A matriz A pode ser expressa como:
 
$$\mathbf{A}_{n \times n} = -\frac{1}{2} \left[ \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}' \right] \mathbf{D}_{n \times n}^2 \left[ \mathbf{I}_n - \frac{1}{n} \mathbf{1} \mathbf{1}' \right]$$
    - $\mathbf{A}_{n \times n}$  é sempre não negativa definida ou positiva definida
    - $\mathbf{D}^2$ : elementos de  $\mathbf{D}$  ao quadrado

Análise Multivariada - 2022

174

- Sejam  $\lambda_1, \lambda_2, \dots, \lambda_n$  os autovalores de  $\mathbf{A}$  e  $\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_n$ , seus correspondentes autovetores
  - ✓ Se  $\mathbf{A}$  tem posto q, então, q autovalores são diferentes de zero e (n – q) são iguais a zero
  - ✓ A matriz A pode ser expressa na forma:
 
$$\mathbf{A}_{n \times n} = \mathbf{V} \mathbf{\Lambda} \mathbf{V}' \quad \mathbf{V} = [\mathbf{v}_1 \quad \mathbf{v}_2 \quad \dots \quad \mathbf{v}_n]$$
  - ✓ Considerando-se os q autovalores não-nulos, definem-se as novas coordenadas (estímulos):
 
$$\mathbf{Y}_j = \mathbf{v}_j \sqrt{\lambda_j}$$
    - A matriz das novas coordenadas é:  $\mathbf{Y}_{n \times q} = \mathbf{V}_{n \times q}^* \mathbf{\Lambda}_{q \times q}^{1/2}$
  - ✓ A matriz de dados original  $\mathbf{X}_{n \times p}$  é transformada em  $\mathbf{Y}_{n \times q}$ .

Análise Multivariada - 2022

175

- Distâncias Euclidianas originais:

$$d_{ij}^2 = (\mathbf{X}_i - \mathbf{X}_j)(\mathbf{X}_i - \mathbf{X}_j)' = (\mathbf{Y}_i - \mathbf{Y}_j)(\mathbf{Y}_i - \mathbf{Y}_j)'$$

✓ Se ao invés de construir q novas coordenadas, forem construídas r,  $r < q$ , as novas distâncias não serão exatamente iguais às originais

- Poderão ser tratadas como uma aproximação
- Necessária avaliação da qualidade da aproximação

Análise Multivariada - 2022

176

### Exemplo

- Matriz de distâncias:

$$\mathbf{D}_{3 \times 3} = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 0 & 3 \\ 2 & 3 & 0 \end{bmatrix} \cdot \mathbf{D}_{3 \times 3}^2 = \begin{bmatrix} 0 & 1 & 4 \\ 1 & 0 & 9 \\ 4 & 9 & 0 \end{bmatrix}.$$

- Matriz  $\mathbf{A}$ :

$$a_{ij} = -\frac{1}{2} [d_{ij}^2 - d_{i.}^2 - d_{.j}^2 + d_{..}^2] \cdot \mathbf{A}_{3 \times 3} = \begin{bmatrix} 0,111 & 0,444 & -0,556 \\ 0,444 & 1,778 & -2,222 \\ -0,556 & -2,222 & 2,778 \end{bmatrix}.$$

✓ Autovalores de  $\mathbf{A}$ :  $\lambda_1 = 4,667$  e  $\lambda_2 = \lambda_3 = 0$

✓ Obtém-se uma nova coordenada ( $\mathbf{Y}_1$ )

Análise Multivariada - 2022

177

- Matrizes  $\mathbf{V}$  e  $\mathbf{V}^*$ :

$$\mathbf{V}_{3 \times 3} = \begin{bmatrix} -0,154 & -0,407 & 0,900 \\ 0,617 & -0,672 & -0,409 \\ 0,772 & -0,618 & -0,147 \end{bmatrix} \cdot \mathbf{V}_{3 \times 1}^* = \begin{bmatrix} -0,154 \\ 0,617 \\ 0,772 \end{bmatrix}.$$

✓ Vetor com novas coordenadas:

$$\mathbf{Y}_1 = \mathbf{v}_1 \sqrt{\lambda_1} = \begin{bmatrix} -0,154 \\ 0,617 \\ 0,772 \end{bmatrix} \sqrt{4,667} = \begin{bmatrix} -0,333 \\ -1,333 \\ 1,667 \end{bmatrix}.$$

– Informação do 1º vetor observado está contida em  $-0,333$

✓ Distâncias originais foram preservadas:

$$d_{12}^2 = (-0,333 + 1,333)^2 = 1$$

$$d_{13}^2 = (-0,333 - 1,667)^2 = 4$$

$$d_{23}^2 = (-1,333 - 1,667)^2 = 9.$$

Análise Multivariada - 2022

178

### Método MDS – Escalonamento Não-métrico

- Há um total de  $M = n(n - 1)/2$  dissimilaridades entre pares de n objetos diferentes

✓ Considere as dissimilaridades ordenadas

$$\delta_{i_1, k_1} < \delta_{i_2, k_2} < \dots < \delta_{i_M, k_M}.$$

– Sem perda de generalidade, suponha que não haja empates

- Objetivo:

✓ Encontrar q dimensões que representem a informação de ordenação

Análise Multivariada - 2022

179



• Procedimento:

✓ Função monótona crescente que aproxima as dissimilaridades originais  $\delta_{ij,kj}$  em um novo conjunto de distâncias  $d_{ij,kj}$

$$d_{ij,kj} \approx f(\delta_{ij,kj}).$$

– Onde  $f(\cdot)$  é tal que

$$\delta_{ij,kj} < \delta_{il,kl} \Rightarrow f(\delta_{ij,kj}) < f(\delta_{il,kl}).$$

✓ Algoritmo para escalonamento não métrico :

– Sheppard-Kruskal

Análise Multivariada - 2022

180

## Determinação das Novas Coordenadas

• Para determinados valores de  $q$ , pode-se não encontrar a representação perfeita

• Processo iterativo de ordenação:

✓ Envolve métodos de regressão por mínimos quadrados

✓ Otimiza medida para quantificar a concordância da ordenação de distâncias (dissimilaridades) nas novas dimensões com a ordenação dos valores originais

✓ Minimização por métodos do gradiente e Newton-Raphson

Análise Multivariada - 2022

181

• Stress (Standardized Residual Sum of Squares):

$$\text{Stress}(q) = \left[ \frac{\sum_{i < k} \sum (d_{ik} - d_{ik}^q)^2}{\sum_{i < k} \sum (d_{ik})^2} \right]^{1/2}.$$

✓ Caso métrico

–  $d_{ik}^q$ : distância entre  $i$  e  $k$  nas novas dimensões

–  $d_{ik}$ : distância original entre  $i$  e  $k$

✓ Caso não métrico

– Constantes de referência para implementação do MDS

✓ Busca-se encontrar os novos valores  $d_{ik}^q$  de modo que o coeficiente de Stress seja o menor possível

Análise Multivariada - 2022

182

• Valores de referência – Stress:

Stress (%)	Qualidade de Ajuste
20	Ruim
10	Razoável
5	Bom
2,5	Excelente
0	Perfeito

Análise Multivariada - 2022

183

- Coeficiente SStress de Takame:

$$SStress(q) = \left[ \frac{\sum_{i < k} \sum ((d_{ik})^2 - (d_{ik}^q)^2)^2}{\sum_{i < k} \sum (d_{ik})^4} \right]^{1/2}.$$

- ✓ Em geral, diz-se que o ajuste é bom se  $SStress < 0,1$
- ✓ Algoritmo de implementação:
  - ALSCAL

Análise Multivariada - 2022

184

- Regressão linear:

- ✓ Outra maneira de avaliar a qualidade do ajuste
- ✓ Gráfico de dispersão das distâncias (dissimilaridades) originais vs. as novas
- ✓ Ajuste de modelo de regressão simples
- ✓ Cálculo do coeficiente de determinação ( $R^2$ )
  - Quanto maior o  $R^2$ , melhor o ajuste

Análise Multivariada - 2022

185

### Exemplo 7.3

Mingoti, 2005 – Continuação Ex. 6.8 e 7.1

- Dados relativos a 21 países (ONU, 2002)
- Variáveis:
  - ✓ Expectativa de vida
  - ✓ Educação
  - ✓ Renda (PIB)
  - ✓ Estabilidade política e de segurança
- Método de agrupamento: escalonamento multidimensional.
- Conjunto de dados: *BD\_multivariada.xls/paises*

Análise Multivariada - 2022

186

- Novas coordenadas ( $q = 2$ )

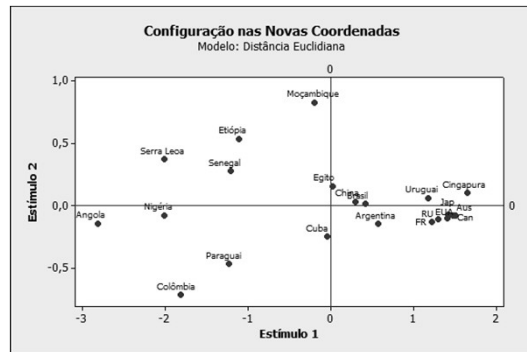
País	Estímulo 1	Estímulo 2
Reino Unido	1,2979	-0,1103
Austrália	1,5022	-0,0766
Canadá	1,4778	-0,082
Estados Uni	1,4045	-0,0996
Japão	1,4221	-0,0745
França	1,2245	-0,1305
Cingapura	1,6504	0,1001
Argentina	0,5706	-0,147
Uruguai	1,1731	0,0624
Cuba	-0,0473	-0,2494
Colômbia	-1,8131	-0,7174
Brasil	0,4138	0,0138
Paraguai	-1,2359	-0,468
Egito	0,0209	0,1569
Nigéria	-2,0142	-0,0764
Senegal	-1,2051	0,2752
Serra Leoa	-2,0144	0,3707
Angola	-2,8117	-0,1437
Etiópia	-1,1136	0,5337
Moçambique	-0,1984	0,8292
China	0,2959	0,0334

Stress: 0,00632  
SStress: 0,00392

Análise Multivariada - 2022

188

• Gráfico do Escalonamento Multidimensional



✓ Solução de agrupamento semelhante às obtidas nos métodos anteriores

Análise Multivariada - 2022

189

• Comentários:

- ✓ Moçambique se mostra um pouco afastado de Senegal, Etiópia e Serra Leoa (mais apropriado grupo isolado)
- ✓ Austrália, Canadá, Cingapura, EUA, Japão, Reino Unido e Uruguai
  - Países muito similares (pontos sobrepostos)
  - Uruguai e Cingapura mais destacados (características um pouco diferentes)
- ✓ Grande similaridade entre Brasil e China
- ✓ Paraguai e Colômbia com características semelhantes
  - Mais distantes de Angola e Nigéria

Análise Multivariada - 2022

191

**Considerações Finais**

**Análise de Agrupamentos**

- É das técnicas mais utilizadas
  - ✓ Nem todos os pacotes dispõem de cálculo automático dos critérios numéricos de ajuste
- Observações incompletas
  - ✓ Possíveis de serem utilizadas
  - ✓ Comparação de modo ponderado dos elementos amostrais

Análise Multivariada - 2022

200

- Útil para redução dimensionalidade da amostra ou do número de variáveis a serem avaliadas
  - ✓ Análise mais elaborada em cada conglomerado é mais informativa do que uma análise do conjunto de  $n$  observações como um todo
  - ✓ Resumo da informação de cada cluster
    - Análise Fatorial, Análise de Componentes Principais

Análise Multivariada - 2022

201

- Alternativa para validação da solução
  - ✓ Amostras aleatórias do conjunto de dados originais
  - ✓ Para cada amostra, aplica-se o método escolhido de agrupamento
  - ✓ Verifica-se se os valores estimados do número de grupos são semelhantes nas várias amostras
- Análise do poder de separação
  - ✓ Análise Discriminante

Análise Multivariada - 2022

202

- Variáveis observadas têm distribuição normal multivariada:
  - ✓ Comparação dos vetores de médias dos grupos formados
    - Testes de significância por MANOVA

Análise Multivariada - 2022

203

## Referências

### **Bibliografia Recomendada**

- MANLY, B. J. F. *Métodos Estatísticos Multivariados: uma Introdução*. Bookman, 2008.
- JOHNSON, R. A.; WINCHERN, D. W. *Applied Multivariate Statistical Analysis*. Prentice Hall, 2007
- MINGOTI, S.A. *Análise de Dados através de Métodos de Estatística Multivariada*. Ed. UFMG, 2005.
- EVERITT, B.; HOTHORN, T. *An Introduction to Applied Multivariate Analysis with R*. Springer, 2011.