

Tópico:

Análise de Conglomerados ou Agrupamentos (ou ainda, Cluster Analysis)

Bibliografia:

- R.A. **Johnson**, Applied Multivariate Statistical Analysis, Prentice Hall, 1992
- L.J. **Corrar**; E. **Paulo**; J.M. **Dias Filho**, Analise multivariada ..., Atlas, 2007
- L.P. **Favero**; et al, Análise de dados - modelagem multivariada ..., Campus, 2009.

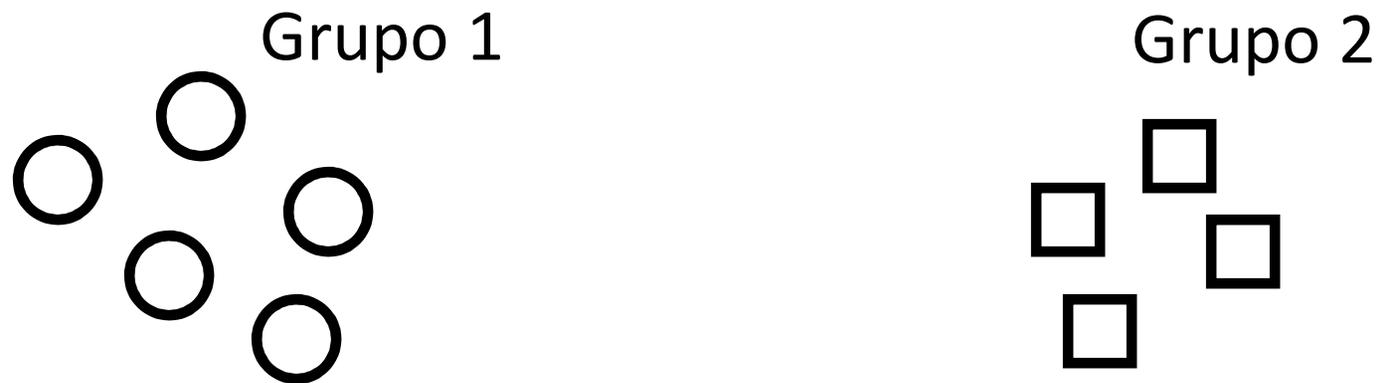
Análise de agrupamentos (cluster): Objetivo dessa técnica é agrupar objetos (itens) semelhantes segundo suas características (variáveis).

Embora vamos considerar **agrupamentos de objetos**, não existem impedimentos para realizar o agrupamento de variáveis semelhantes segundo as realizações obtidas pelos objetos amostrados. Análise fatorial se mostra mais robusta para o agrupamento de variáveis.

Análise de Agrupamentos

Finalidade: Encontrar grupos de modo a

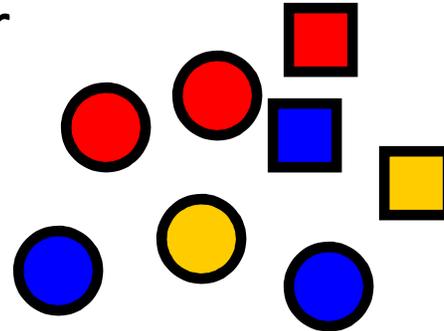
- maximizar a homogeneidade dentro dos grupos
- maximizar a heterogeneidade entre os grupos



Dois grupos classificados pela forma

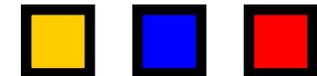
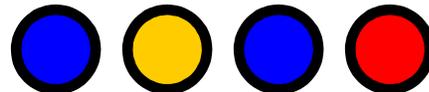
Análise de Agrupamentos

Formar os agrupamentos não é tão simples, pois os objetos são, em geral, definidos por várias características



Possibilidades:

2 grupos:



3 grupos:



A formação dos agrupamentos vai depender da definição de “similaridade”.

- Identificar grupos de investimento de acordo com perfis de risco;
- Identificar segmentos homogêneos de consumidores, a fim de estabelecer programas específicos para cada público;
- Agrupar empresas com base em indicadores financeiros;
- ...
- Exemplo: Häagen-Dazs ([Malhotra 4ª Ed, Cap. 20, p. 571](#))

Análise de Agrupamentos – Qual critério?

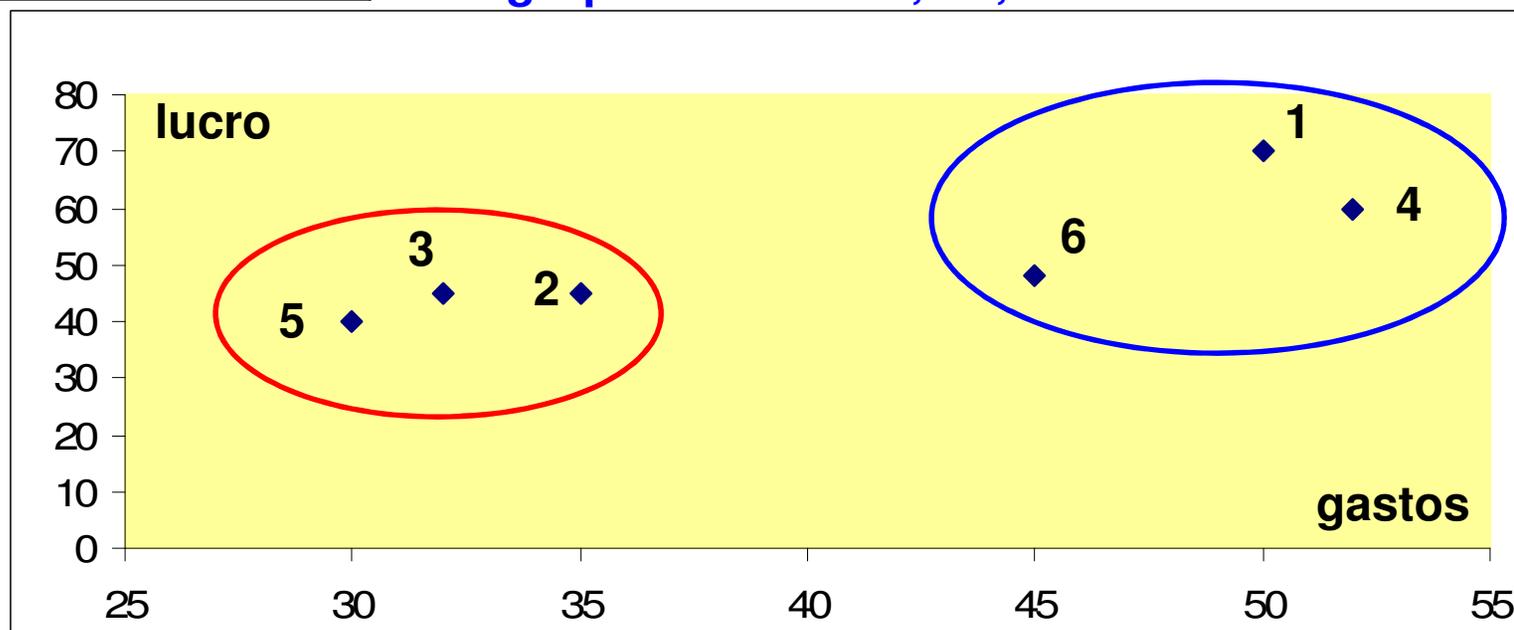
Poderíamos procurar grupos graficamente.

	gastos	lucro
empresa1	50	70
empresa2	35	45
empresa3	32	45
empresa4	52	60
empresa5	30	40
empresa6	45	48

Com duas variáveis poderíamos fazer um diagrama de dispersão e agrupar baseado na distância entre pontos.

Agrupamento 1: E2, E3, E5

Agrupamento 2: E1, E4, E6



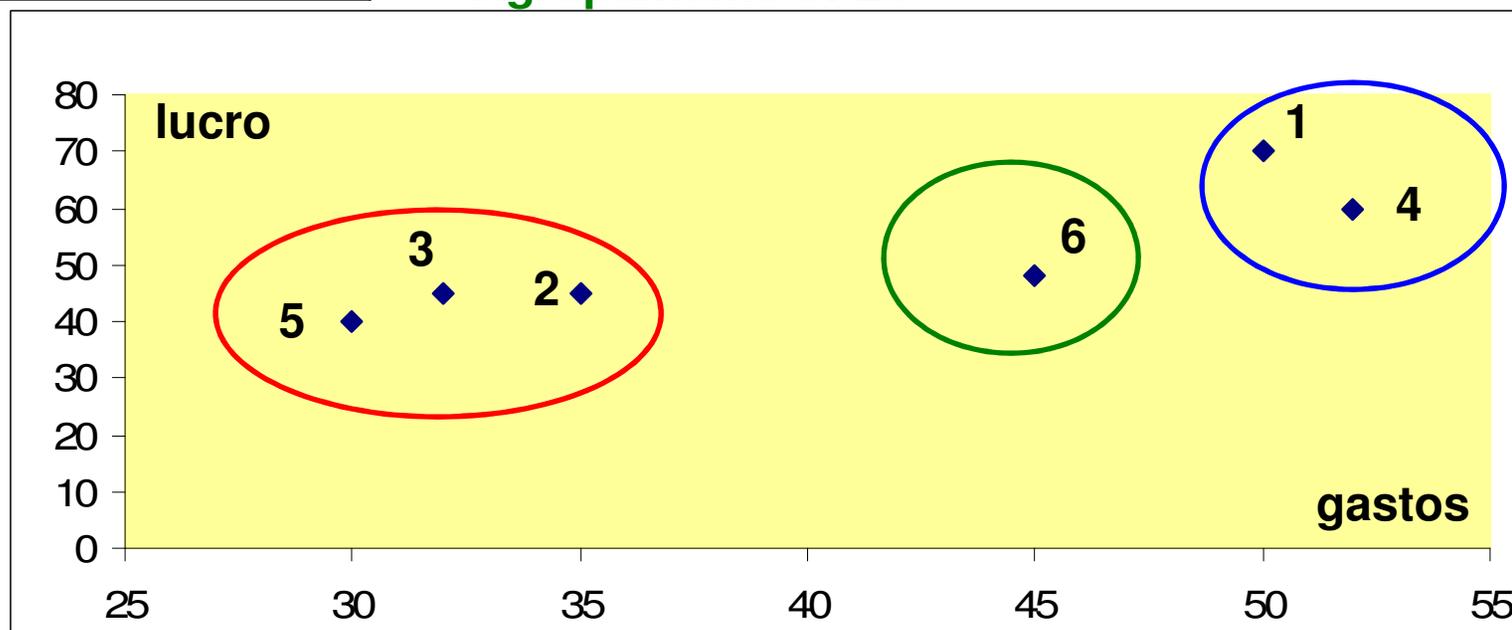
	gastos	lucro
empresa1	50	70
empresa2	35	45
empresa3	32	45
empresa4	52	60
empresa5	30	40
empresa6	45	48

Poderíamos formar três grupos:

Agrupamento 1: E2, E3, E5

Agrupamento 2: E1, E4

Agrupamento 3: E6



Etapas de uma análise de conglomerados

- Análise das variáveis e dos objetos a serem agrupados (seleção das variáveis, identificação de outliers e padronização de variáveis).
- Seleção da medida de **distância** ou **semelhança** entre cada par de objetos.
- Seleção do algoritmo de agrupamento: método **hierárquico** ou **não-hierárquico**.
- Escolha da quantidade de agrupamentos formados interpretação e validação dos agrupamentos.

Variáveis e objetos a serem agrupados

A estrutura que os grupos refletem dependem da seleção das variáveis.

A técnica de análise de agrupamentos não distingue se as variáveis são ou não relevantes (**Essa tarefa fica a cargo do pesquisador**). Sugestões:

- Pesquisas passadas
- Teoria
- Hipóteses que estão sendo testadas
- Presença do julgamento e intuição do pesquisador

A técnica é sensível a observações atípicas (outliers?).

Observações atípicas podem formar grupos isolados, que podem ser interesse do pesquisador.

Padronização das variáveis

- Elimina viés introduzido por diferenças de escala
- Z_scores:
$$Z = \frac{x - \bar{x}}{s}$$
- Recomenda-se quando há diferentes escalas entre as variáveis
- Tem impacto substancial na solução final: padronização minimiza diferenças de grupo, pois se grupos são bem separados por uma dada variável X_i , a variância dessa variável será grande.

Precisamos medir correspondências entre objetos. Ou seja medidas de similaridade e de dissimilaridade.

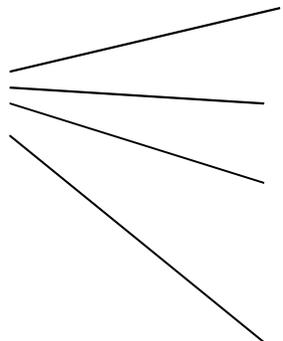
Os métodos mais usados são:

- Medidas de distância (**dissimilaridade**)
- Medidas de correlação (**similaridade**)
- Medidas de associação (**dissimilaridade** ou **similaridade** com variáveis nominais)

Dissimilaridade

Note que estamos interessados em agrupar objetos. As medidas de distância são entre dois objetos (itens).

objetos



	X_1	X_2	...	X_p
Item 1	X_{11}	X_{12}	...	X_{1p}
Item 2	X_{21}	X_{22}	...	X_{2p}
Item 3	X_{31}	X_{32}	...	X_{3p}
...
Item n	X_{n1}	X_{n2}	...	X_{np}

Medidas de distância

Medidas de distância são consideradas medidas de dissimilaridade pois quanto maior a distância, menor a semelhança entre objetos, e vice-versa.

Medidas de distância d_{ij} entre o item i e o item j

Podemos obter d_{ij} a partir das variáveis.

Como?

	X_1	X_2	...	X_p
Item 1	X_{11}	X_{12}	...	X_{1p}
...	X_{21}	X_{22}	...	X_{2p}
Item i	X_{i1}	X_{i2}	...	X_{ip}
...
Item j	X_{j1}	X_{j2}	...	X_{jp}
...
Item n	X_{n1}	X_{n2}	...	X_{np}

Algumas medidas de distância

$$\begin{aligned} \text{Distância euclidiana : } d_{ij} &= \sqrt{\sum_{k=1}^p (x_{ik} - x_{jk})^2} \\ &= \sqrt{(x_{i1} - x_{j1})^2 + (x_{i2} - x_{j2})^2 + \dots + (x_{ip} - x_{jp})^2} \end{aligned}$$

$$\text{Distância quadrática euclidiana : } d_{ij}^2 = \sum_{k=1}^p (x_{ik} - x_{jk})^2$$

$$\text{Distância euclidiana média : } d_{ij} = \sqrt{\frac{\sum_{k=1}^p (x_{ik} - x_{jk})^2}{p}}$$

$$\text{Distância euclidiana padronizada : } d_{ij} = \sqrt{\sum_{k=1}^p \left(\frac{x_{ik} - x_{jk}}{\sqrt{s_{kk}}} \right)^2}$$

Algumas medidas de distância

Distância de Minkowski:
$$d_{ij} = \left(\sum_{k=1}^p |x_{ik} - x_{jk}|^n \right)^{\frac{1}{n}}$$

$n = 1$: é conhecida por “métrica do quarteirão” (métrica city-block)

$n = 2$: representa a distância euclidiana

variações de n causam trocas nos pesos dados a pequenas e a grandes diferenças.

Distância de Mahalanobis:
$$d_{ij} = \sqrt{(x_i - x_j)' S^{-1} (x_i - x_j)}$$

onde: $x'_i = [x_{i1}, x_{i2}, x_{i3}, \dots, x_{ip}]$

$x'_j = [x_{j1}, x_{j2}, x_{j3}, \dots, x_{jp}]$

S^{-1} = inversa da matriz de covariância amostral

$$\text{Correlação de Pearson : } r_{ij} = \frac{\sum_{k=1}^p (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j)}{\sqrt{\sum_{k=1}^p (x_{ik} - \bar{x}_i)^2 \sum_{k=1}^p (x_{jk} - \bar{x}_j)^2}}$$

onde: x_{ik} = valor da variável k para o objeto i

x_{jk} = valor da variável k para o objeto j

\bar{x}_i = média das variáveis para o objeto i

\bar{x}_j = média das variáveis para o objeto j

p = número de variáveis

O valor de r varia entre -1 e 1, sendo que o zero significa que não há associação. Quanto maiores as correlações, mais próximas as observações.

Exemplo

Um exemplo: vamos encontrar a **matriz de distâncias** para o seguinte exemplo:

	gasto	lucro
Empresa 1	50	70
Empresa 2	35	45
Empresa 3	32	45

Vamos utilizar a distância quadrática:

$$d_{ij}^2 = \sum_{k=1}^p (x_{ik} - x_{jk})^2$$

Temos:

distância entre empresa 1 e empresa 2 : $d_{12}^2 = (50 - 35)^2 + (70 - 45)^2 = 850$

distância entre empresa 1 e empresa 3 : $d_{13}^2 = (50 - 32)^2 + (70 - 45)^2 = 949$

distância entre empresa 2 e empresa 3 : $d_{23}^2 = (35 - 32)^2 + (45 - 45)^2 = 9$

Matriz de distância para o exemplo

$$D = \begin{pmatrix} d_{11} & d_{12} & d_{13} \\ d_{21} & d_{22} & d_{23} \\ d_{31} & d_{32} & d_{33} \end{pmatrix} = \begin{pmatrix} 0 & 850 & 949 \\ 850 & 0 & 9 \\ 949 & 9 & 0 \end{pmatrix}$$

Propriedades observadas no exemplo:

-elementos nulos na diagonal (distância entre um mesmo objeto é nula)

-Matriz simétrica (distância de i até j = distância de j até i)

Os agrupamentos hierárquicos: são realizados por sucessivas fusões ou por sucessivas divisões.

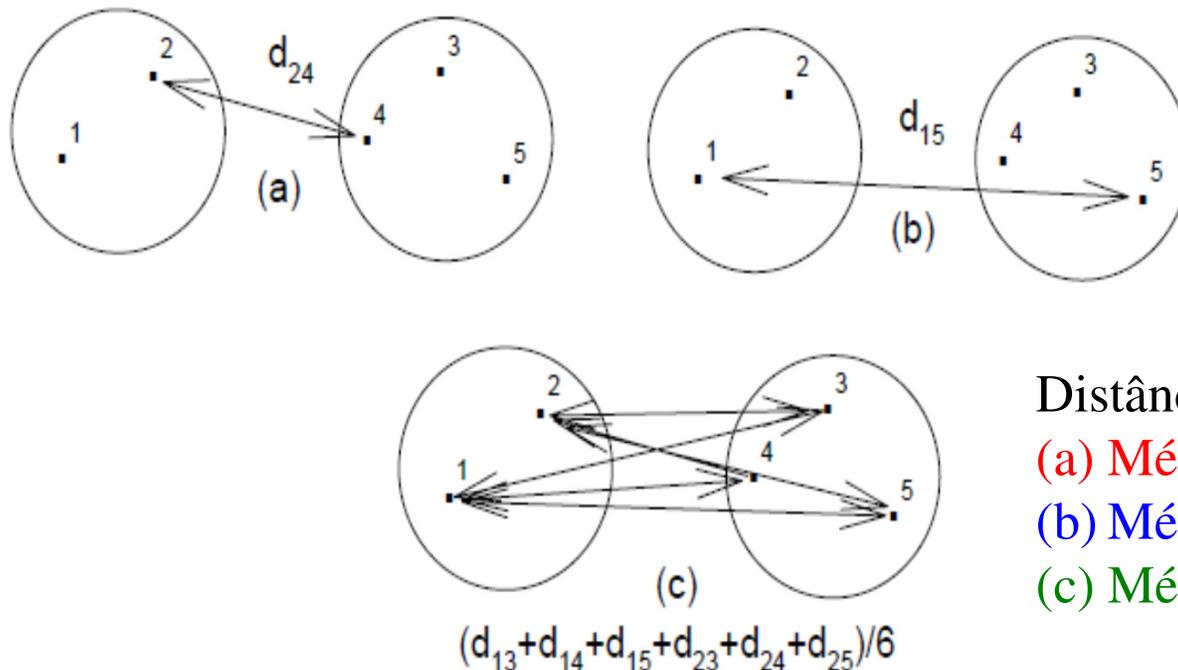
métodos hierárquicos aglomerativos: iniciam com tantos grupos quanto aos objetos, ou seja, cada objeto forma um agrupamento. Inicialmente, os objetos mais similares são agrupados e fundidos formando um único grupo. Eventualmente o processo é repetido, e com o decréscimo da similaridade, todos os subgrupos são fundidos, formando um único grupo com todos os objetos.

Os métodos hierárquicos divisivos: Um único subgrupo inicial existe com todos os objetos e estes são subdivididos em dois subgrupos de tal forma que exista o máximo de semelhança entre os objetos dos mesmos subgrupos e a máxima dissimilaridade entre elementos de subgrupos distintos. Estes subgrupos são posteriormente subdivididos em outros subgrupos dissimilares. O processo é repetido até que haja tantos subgrupos quantos objetos.

Consideraremos agrupamentos **hierárquicos aglomerativos** (“Linkage Methods”)

Vamos considerar:

- Método de ligação simples (mínima distância ou vizinho mais próximo),
- Método ligação completa (máxima distância ou vizinho mais distante) e
- Método da ligação média (distância média).



Distâncias entre os grupos para:
(a) Método da ligação simples,
(b) Método da ligação completa e
(c) Método da ligação média.

Algoritmo para grupamentos hierárquicos aglomerativos com n objetos (itens ou variáveis).

1. Iniciar com n grupos, cada um com um único elemento e com uma matriz simétrica $n \times n$ de dissimilaridades (distâncias) $D = \{d_{U,V}\}$.
2. Buscar na matriz D o par de grupos mais similar (menor distância) e fazer a distância entre os grupos mais similares U e V igual a $d_{U,V}$.
3. Fundir os grupos U e V e nomeá-lo por (UV) . Recalcular e rearranjar as distâncias na matriz D (a) eliminando as linhas e colunas correspondentes a U e V e (b) acrescentando uma linha e coluna com as distâncias entre o grupo (UV) e os demais grupos.
4. Repetir os passos 2 e 3 num total de $(n-1)$ vezes (todos os objetos estarão em único grupo). Anotar a identidade dos grupos que vão sendo fundidos e os respectivos níveis (distâncias) nas quais isto ocorre.

No passo 3 será usado: $d_{(UV),W} = \min\{d_{U,W}, d_{V,W}\}$.

Vamos utilizar um exemplo onde temos 4 objetos (A, B, C, D) com a seguinte matriz de distância:

$$D = \begin{matrix} & \begin{matrix} A & B & C & D \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} 0 & & & \\ 3 & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{bmatrix} \end{matrix}$$

Os objetos menos distantes devem, inicialmente, ser agrupados.

Mínima distância entre i e $j = \min(d_{i,j}) = d_{A,B} = 3$.

Próximo passo: juntar A com B formando o grupo (AB) e em seguida calcular as distâncias deste grupo e os objetos remanescentes.

As distâncias dos vizinhos mais próximos são:

$$d_{(AB),C} = \min\{d_{A,C}, d_{B,C}\} = \min\{7, 9\} = 7$$

$$d_{(AB),D} = \min\{d_{A,D}, d_{B,D}\} = \min\{8, 6\} = 6$$

$$D = \begin{matrix} & A & B & C & D \\ A & \begin{bmatrix} 0 & & & \end{bmatrix} \\ B & \begin{bmatrix} 3 & 0 & & \end{bmatrix} \\ C & \begin{bmatrix} 7 & 9 & 0 & \end{bmatrix} \\ D & \begin{bmatrix} 8 & 6 & 5 & 0 \end{bmatrix} \end{matrix}$$

E obtemos a nova matriz D para o próximo passo:

$$D = \begin{matrix} & AB & C & D \\ AB & \begin{bmatrix} 0 & & \end{bmatrix} \\ C & \begin{bmatrix} 7 & 0 & \end{bmatrix} \\ D & \begin{bmatrix} 6 & 5 & 0 \end{bmatrix} \end{matrix}$$

A menor distância é entre D e C, $d_{D,C} = 5$. Então DC formam um grupo no nível (distância) 5.

Recalculando a distância:

$$d_{(DC),(AB)} = \min\{d_{D,(AB)}, d_{C,(AB)}\}$$

$$= \min\{6, 7\} = 6$$

$$D = \begin{array}{c} \begin{array}{cc} & \begin{array}{ccc} AB & C & D \end{array} \\ \begin{array}{c} AB \\ C \\ D \end{array} & \begin{bmatrix} 0 & & \\ 7 & 0 & \\ 6 & 5 & 0 \end{bmatrix} \end{array}$$

E a nova matriz D será:

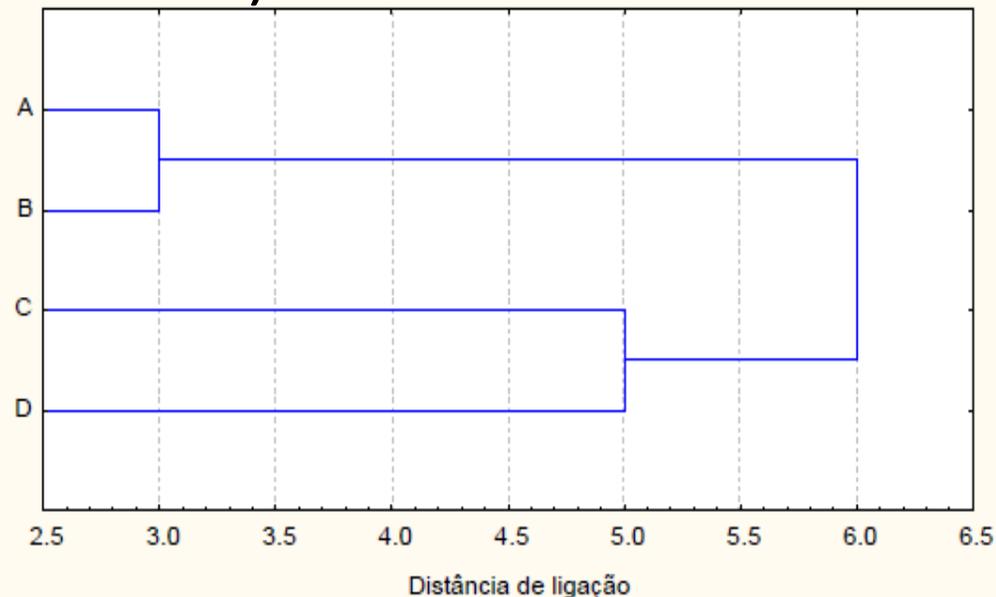
$$D = \begin{array}{c} \begin{array}{cc} & \begin{array}{cc} DC & AB \end{array} \\ \begin{array}{c} DC \\ AB \end{array} & \begin{bmatrix} 0 & \\ 6 & 0 \end{bmatrix} \end{array}$$

O grupo DC é agrupado com o grupo AB na distância 6

Dendrograma (ou diagrama de árvore):

Os agrupamentos são mostrados graficamente. Os ramos representam os agrupamentos e os nodos representam os pontos onde os agrupamentos se fundem ao longo do eixo de distâncias (ou similaridade)

No exemplo visto:



Dendrograma para agrupar 4 objetos (A, B, C e D) pelo método da ligação simples (vizinho mais próximo).

Parecido com o anterior, mas a distância entre grupos é tomada como a “máxima” distância entre dois elementos de cada grupo.

No passo 3 é usado: $d_{(UV),W} = \max\{d_{U,W}, d_{V,W}\}$

Vamos utilizar o mesmo exemplo com 4 objetos (A, B, C, D) e matriz de distância:

Os objetos menos distantes devem, inicialmente, ser agrupados.

Mínima distância entre i e j

$$= \min(d_{i,j}) = d_{A,B} = 3.$$

Próximo passo: juntar A com B formando o grupo (AB) e em seguida calcular as distâncias deste grupo e os objetos remanescentes.

$$D = \begin{matrix} & A & B & C & D \\ A & [0 & & &] \\ B & 3 & 0 & & \\ C & 7 & 9 & 0 & \\ D & 8 & 6 & 5 & 0 \end{matrix}$$

As distâncias dos vizinhos mais distantes são:

$$d_{(AB),C} = \max\{d_{A,C}, d_{B,C}\} = \max\{7, 9\} = 9$$

$$d_{(AB),D} = \max\{d_{A,D}, d_{B,D}\} = \max\{8, 6\} = 8$$

$$D = \begin{matrix} & A & B & C & D \\ A & \left[\begin{array}{cccc} 0 & & & \\ 3 & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{array} \right] \\ B & & & & \\ C & & & & \\ D & & & & \end{matrix}$$

E obtemos a nova matriz D para o próximo passo:

A menor distância é entre D e C, $d_{D,C} = 5$.

Então DC formam um grupo no nível 5.

Recalculando a distância:

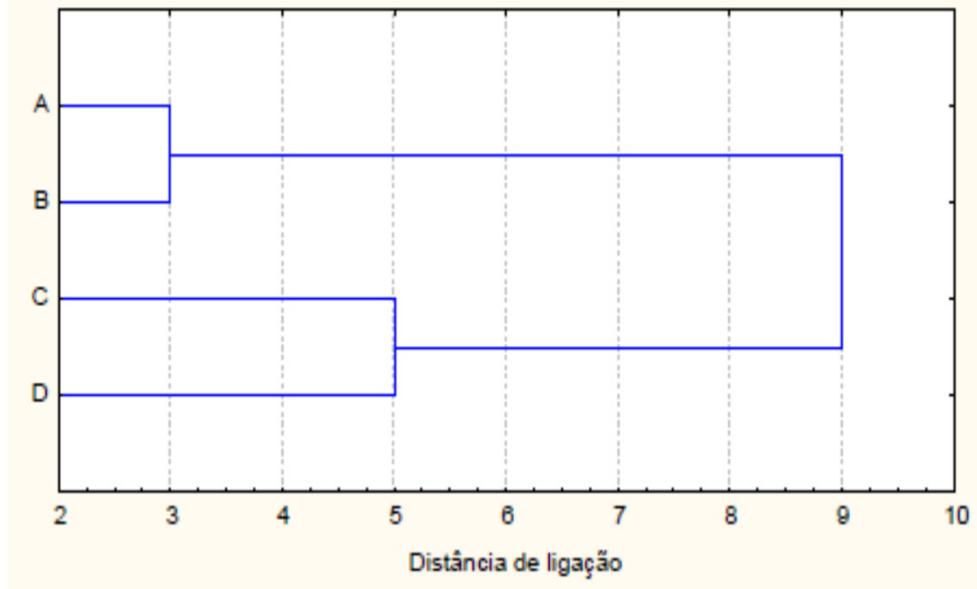
$$\begin{aligned} d_{(DC),(AB)} &= \max\{d_{D,(AB)}, d_{C,(AB)}\} \\ &= \min\{8, 9\} = 9 \end{aligned}$$

$$D = \begin{matrix} & AB & C & D \\ AB & \left[\begin{array}{ccc} 0 & & \\ 9 & 0 & \\ 8 & 5 & 0 \end{array} \right] \\ C & & & \\ D & & & \end{matrix}$$

E a nova matriz D será:

$$D = \begin{matrix} & \text{DC} & \text{AB} \\ \text{DC} & \begin{bmatrix} 0 & \\ & \end{bmatrix} \\ \text{AB} & \begin{bmatrix} 9 & 0 \end{bmatrix} \end{matrix}$$

O grupo DC é agrupado com o grupo AB a distância 9



Dendrograma para agrupar 4 objetos (A, B, C e D) pelo método da ligação completa (vizinho mais distante).

Análogo aos anteriores, com exceção de que a distância entre grupos é tomada como a média da distância entre dois elementos de cada grupo.

No passo 3 é usado: $d_{(UV),W} = \text{média}\{d_{U,W}, d_{V,W}\}$

Vamos utilizar o mesmo exemplo com 4 objetos (A, B, C, D) e matriz de distância:

Os objetos menos distantes são agrupados.

Mínima distância entre i e j

$$= \min(d_{i,j}) = d_{A,B} = 3.$$

$$D = \begin{matrix} & \begin{matrix} A & B & C & D \end{matrix} \\ \begin{matrix} A \\ B \\ C \\ D \end{matrix} & \begin{bmatrix} 0 & & & \\ 3 & 0 & & \\ 7 & 9 & 0 & \\ 8 & 6 & 5 & 0 \end{bmatrix} \end{matrix}$$

Próximo passo: juntar A com B formando o grupo (AB) e em seguida calcular as distâncias deste grupo e os objetos remanescentes.

As distâncias entre os grupos são baseadas na média das distâncias entre todos os elementos de um grupo com relação ao elementos de outro grupo

$$d_{(AB),C} = (d_{A,C} + d_{B,C})/2 = (7 + 9)/2 = 8$$

$$d_{(AB),D} = (d_{A,D} + d_{B,D})/2 = (8+6)/2 = 7$$

$$D = \begin{matrix} & A & B & C & D \\ A & 0 & & & \\ B & 3 & 0 & & \\ C & 7 & 9 & 0 & \\ D & 8 & 6 & 5 & 0 \end{matrix}$$

E obtemos a nova matriz D para o próximo passo:

A menor distância é entre D e C, $d_{D,C} = 5$.

Então DC formam um grupo no nível 5.

$$D = \begin{matrix} & AB & C & D \\ AB & 0 & & \\ C & 8 & 0 & \\ D & 7 & 5 & 0 \end{matrix}$$

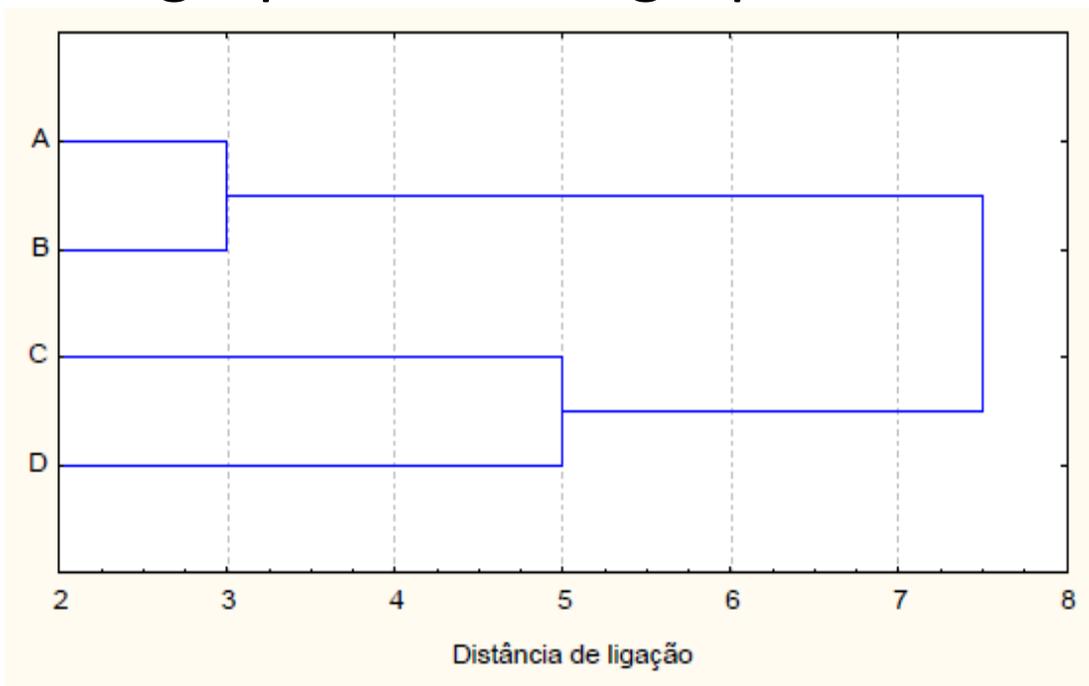
Recalculando a distância:

$$d_{(DC),(AB)} = (d_{D,(AB)} + d_{C,(AB)})/2 = (7+8)/2 = 7,5$$

E a nova matriz D será:

$$D = \begin{matrix} & DC & AB \\ DC & \begin{bmatrix} 0 & \\ & \end{bmatrix} \\ AB & \begin{bmatrix} 7,5 & 0 \end{bmatrix} \end{matrix}$$

O grupo DC é agrupado com o grupo AB a distância 7,5



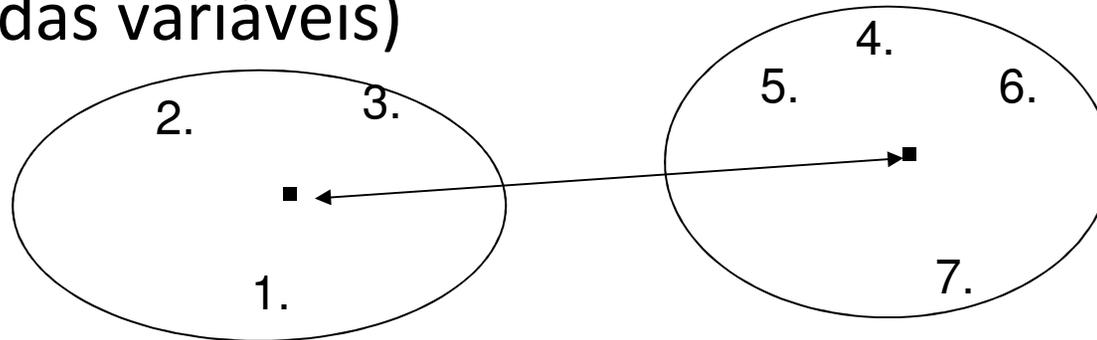
Dendrograma para agrupar 4 objetos (A, B, C e D) pelo método da ligação média.

Note que os métodos podem produzir:

- Agrupamentos com objetos alocados de maneiras diferentes.
- Distâncias de fusão entre grupos diferentes.

Para o exemplo anteriormente considerado, os três métodos não produziram diferenças na alocação dos objetos mas apenas nas distâncias de fusão dos grupos.

Método do centróide: baseia-se na distância entre centróides (média das variáveis)



Método de Ward: Tende a produzir agrupamentos com aproximadamente mesmo número de itens.

- São calculadas médias (centróides) para cada grupo.
- É calculado o quadrado das distâncias entre as médias e os valores das variáveis para cada item.
- Somam-se os quadrados: variância
- Minimiza-se a variância dentro dos grupos: Error Sum of Squares (ESS).

Métodos não hierárquicos

- Formar agrupamentos cujo número de grupos é definido pelo pesquisador, ou seja, particionar n objetos em k grupos.
- Busca coesão interna (semelhança) e isolamento (separação) dos grupos formados.
- Iniciam-se : (1) partição inicial dos itens em grupos ou (2) uma semente inicial de pontos, que formarão os núcleos dos agrupamentos

Algoritmo não hierárquico: K-means

Um dos mais populares algoritmos K-médias (K-means): aloca objetos ao grupo com centróide mais próximo.

Procedimento:

- (1) Particione os itens em k grupos iniciais arbitrariamente;
- (2) Percorra a lista de itens e calcule as distâncias (usualmente a euclidiana) de cada um deles para o centróide (médias) dos grupos. Faça a realocação do item para o grupo em que ele apresentar mínima distância, obviamente se não for o grupo ao qual este pertença. Recalcule os centróides dos grupos que ganharam e perderam item.
- (3) Repita o passo 2 até que nenhuma alteração seja feita.

Algoritmo não hierárquico: K-means

Exemplo:

Utilizando 4 itens (A, B, C e D) e 2 variáveis (X1 e X2) dividir em $k=2$ grupos, pelo método das k-médias.

Objeto	Observação	
	X ₁	X ₂
A	2	0
B	5	2
C	1	4
D	8	4

(1) particionar os itens arbitrariamente em 2 grupos, como por exemplo AD e BC. Calcular a média (centróide) de cada grupo.

Objeto	Centróide	
	\bar{X}_1	\bar{X}_2
AD	$(2+8)/2=5$	$(0+4)/2=2$
BC	$(1+5)/2=3$	$(2+4)/2=3$

Algoritmo não hierárquico: K-means

Recalculando as distâncias dos objetos para o centróide dos grupos e checando a possibilidade de realocação, tem-se:

$$d_{A,D}^2 = 52$$

$$d_{B,D}^2 = 13$$

$$d_{C,D}^2 = 49$$

$$d_{A,(ABC)}^2 = 4,44$$

$$d_{B,(ABC)}^2 = 5,44$$

$$d_{C,(ABC)}^2 = 6,77$$

Nenhuma realocação deve ser realizada, pois os objetos têm menor distância para os respectivos grupos aos quais eles pertencem.

Grupo	Item (distância quadrática p/ centróide)			
	A	B	C	D
D	52,0	13,0	49,0	0,0
ABC	4,4	5,4	6,8	32,4

Para realizar uma checagem da estabilidade de a partição alcançada é recomendável executar novamente o algoritmo com uma nova partição inicial.

- Não há uma resposta certa: depende da situação
- Hierárquico exige menos capacidade de processamento e evidencia *outliers*
- Não-hierárquicos exigem escolhas de sementes, mas são menos suscetíveis a *outliers*, à escolha de medidas e à variáveis inapropriadas
- Tomar um conjunto de dados com uma estrutura de grupos conhecida e ver se o algoritmo é capaz de reproduzir essa estrutura.

**Escolha da quantidade de agrupamentos formados
interpretação e validação dos agrupamentos.**

Não existe um procedimento padrão.

Sugestões:

- Olhar várias medidas de similaridade
- Fazer julgamento teórico
- Fazer análise para diferentes amostras
- Utilizar novas variáveis que teoricamente comprovam os clusters