

Agrupamento de Dados e Aplicações

Algoritmos Hierárquicos (Parte I)

Prof. Eduardo Raul Hruschka



Créditos

- O material a seguir consiste de adaptações e extensões dos originais:
 - Elaborados por Eduardo Hruschka e Ricardo Campello
 - de (Tan et al., 2006)
 - de E. Keogh (SBBD 2003)
- Algumas figuras foram gentilmente cedidas por Lucas Vendramin



Agenda

- Algoritmos Hierárquicos
 - Conceitos e Definições
 - Dendrogramas
 - Grafos de Proximidade
 - *Cophenetic matrices*
 - Métodos Aglomerativos
 - *Single Linkage*
 - *Complete Linkage*
 - Relação com Teoria dos Grafos

Relembrando...

- **Matriz de Dados \mathbf{X} :**

- N linhas (objetos) e n colunas (atributos):

$$\mathbf{X} = \begin{bmatrix} x_{11} & x_{12} & \cdots & x_{1n} \\ x_{21} & x_{22} & \cdots & x_{2n} \\ \vdots & & \ddots & \vdots \\ x_{N1} & x_{N2} & \cdots & x_{Nn} \end{bmatrix}$$

- Cada **objeto** (linha da matriz) é denotado por um vetor \mathbf{x}_i

- Exemplo:

$$\mathbf{x}_1 = [x_{11} \quad \cdots \quad x_{1n}]$$

Relembrando...

- **Matriz de Proximidade** (Dissimilaridade ou Similaridade):

- N linhas e N colunas:

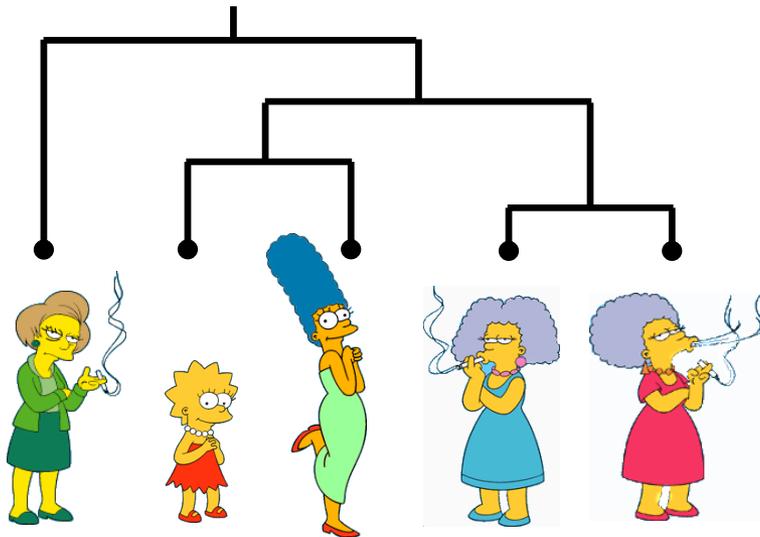
$$\mathbf{D} = \begin{bmatrix} d(\mathbf{x}_1, \mathbf{x}_1) & d(\mathbf{x}_1, \mathbf{x}_2) & \cdots & d(\mathbf{x}_1, \mathbf{x}_N) \\ d(\mathbf{x}_2, \mathbf{x}_1) & d(\mathbf{x}_2, \mathbf{x}_2) & \cdots & d(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & & \ddots & \vdots \\ d(\mathbf{x}_N, \mathbf{x}_1) & d(\mathbf{x}_N, \mathbf{x}_2) & \cdots & d(\mathbf{x}_N, \mathbf{x}_N) \end{bmatrix}$$

- Simétrica se proximidade d apresentar propriedade de simetria

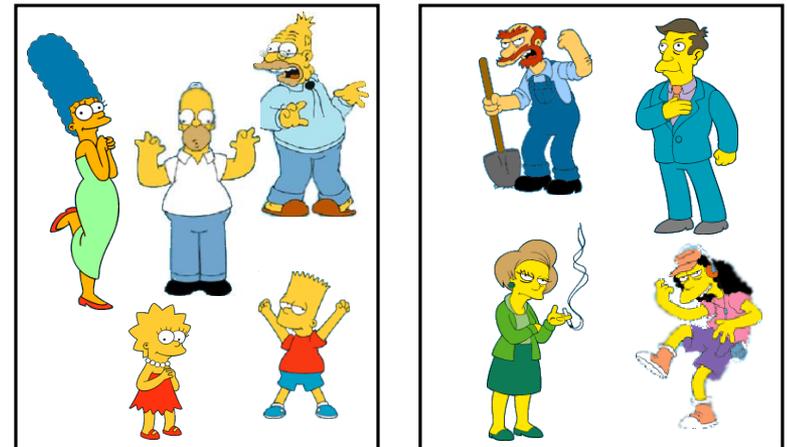
Relembrando...

- **Agrupamento Particional:** constrói uma *partição* dos dados
- **Agrupamento Hierárquico:** constrói uma *hierarquia de partições*

Hierárquicos



Particionais



Definição de Partição de Dados

- Consideremos um conjunto de N objetos a serem agrupados: $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$
- **Partição** (rígida): coleção de k grupos não sobrepostos $\mathbf{P} = \{\mathbf{C}_1, \mathbf{C}_2, \dots, \mathbf{C}_k\}$ tal que:

$$\mathbf{C}_1 \cup \mathbf{C}_2 \cup \dots \cup \mathbf{C}_k = \mathbf{X}$$

$$\mathbf{C}_i \neq \emptyset$$

$$\mathbf{C}_i \cap \mathbf{C}_j = \emptyset \text{ para } i \neq j$$

- Exemplo: $\mathbf{P} = \{ (\mathbf{x}_1), (\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$

Definição de Hierarquia

➤ Hierarquia (de partições de dados):

➤ Sequência de partições aninhadas

- Uma partição \mathbf{P}_1 está *aninhada* em \mathbf{P}_2 se cada grupo de \mathbf{P}_1 é um subconjunto de um grupo de \mathbf{P}_2

➤ Exemplo:

$$\mathbf{P}_1 = \{ (\mathbf{x}_1), (\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$$

$$\mathbf{P}_2 = \{ (\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$$

➤ Contra-Exemplo:

$$\mathbf{P}_3 = \{ (\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$$

$$\mathbf{P}_4 = \{ (\mathbf{x}_1, \mathbf{x}_2), (\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_5) \}$$

Definição de Hierarquia

- Uma hierarquia completa:
 - Inicia ou termina com partição totalmente disjunta
 - *Disjoint clustering*: apenas grupos atômicos (*singletons*)
 - Exemplo: $\mathbf{P} = \{ (\mathbf{x}_1), (\mathbf{x}_2), (\mathbf{x}_3), (\mathbf{x}_4), (\mathbf{x}_5), (\mathbf{x}_6) \}$
 - Também denominada “solução trivial”
 - Inicia ou termina com partição totalmente conjunta
 - *Conjoint clustering*: grupo único com todos os objetos
 - Exemplo: $\mathbf{P} = \{ (\mathbf{x}_1, \mathbf{x}_2, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_5, \mathbf{x}_6) \}$
 - Geralmente possui $N - 2$ partições intermediárias (homework)

Hierarquias podem ser usadas para organizar informação, como, por exemplo, num portal

Web Site Directory - Sites organized by subject

[Suggest your site](#)

Business & Economy

[B2B](#), [Finance](#), [Shopping](#), [Jobs](#)...

Regional

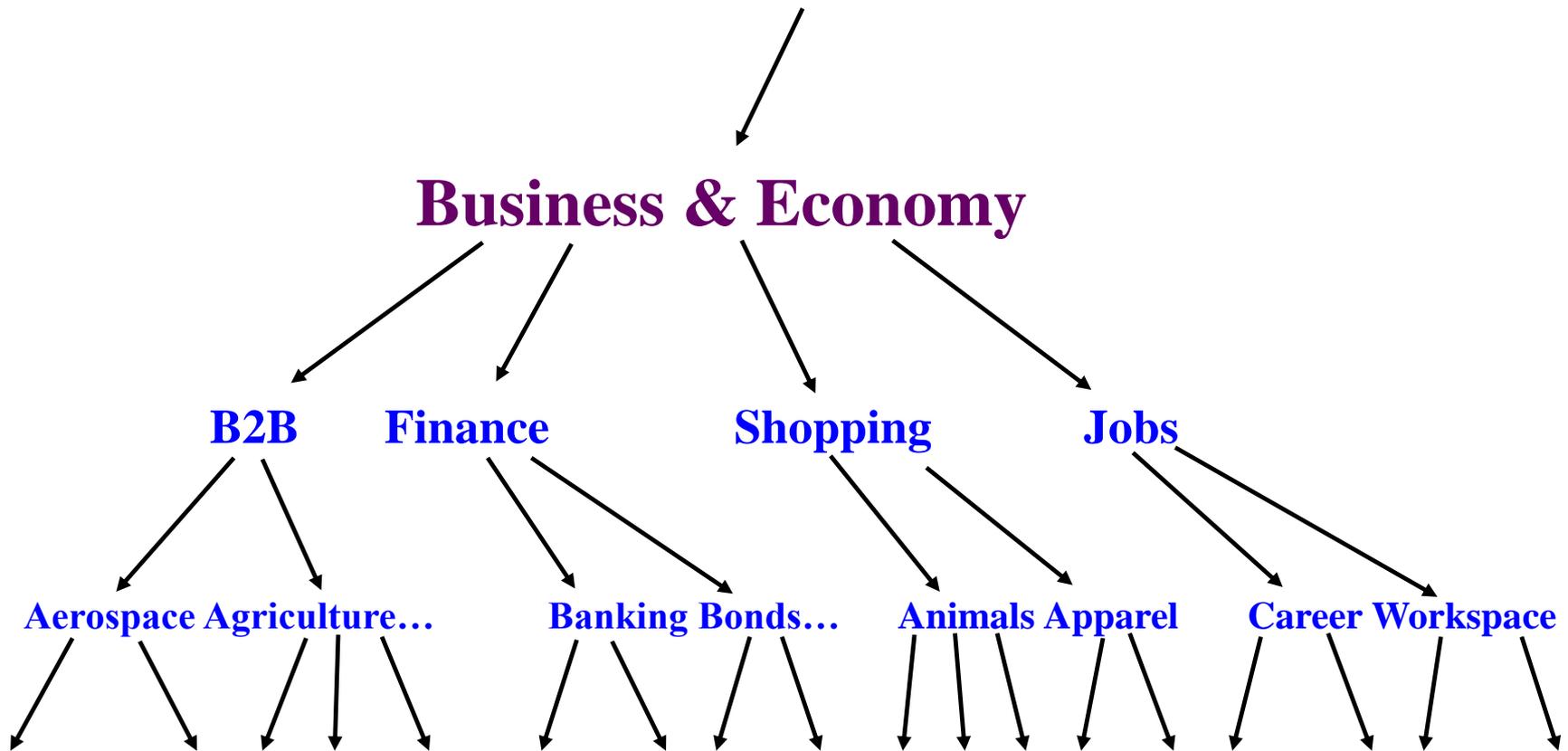
[Countries](#), [Regions](#), [US States](#)...

Computers & Internet

[Internet](#), [WWW](#), [Software](#), [Games](#)...

Society & Culture

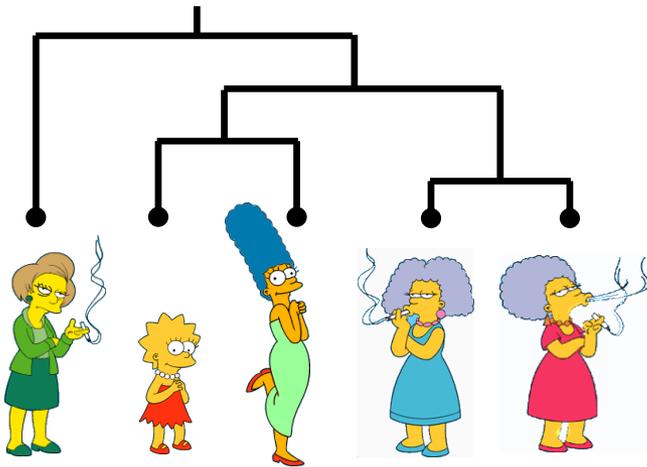
[People](#), [Environment](#), [Religion](#)...



Métodos Clássicos para Agrupamento Hierárquico

Bottom-Up (aglomerativos):

- Iniciar colocando cada objeto em um *cluster*
- Encontrar o melhor par de *clusters* e uni-los
- Repetir até que todos os objetos estejam reunidos num só *cluster*



Top-Down (divisivos):

- Iniciar com todos os objetos em um único *cluster*
- Sub-dividir o *cluster* em dois novos *clusters*
- Aplicar o algoritmo recursivamente em ambos, até que cada objeto forme um *cluster* por si só

Algoritmos hierárquicos
podem operar somente sobre
uma matriz de distâncias: são
relacionais.

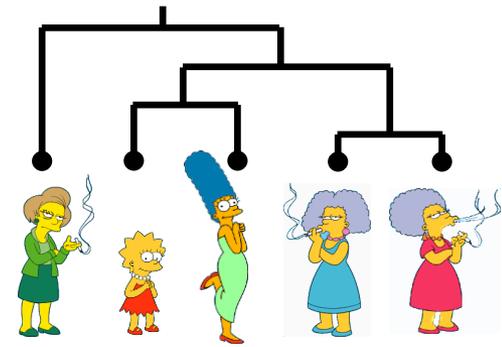
$$D(\text{Mrs. Muntz}, \text{Lisa Simpson}) = 8$$

$$D(\text{Mrs. Krabappel}, \text{Mrs. Krabappel}) = 1$$

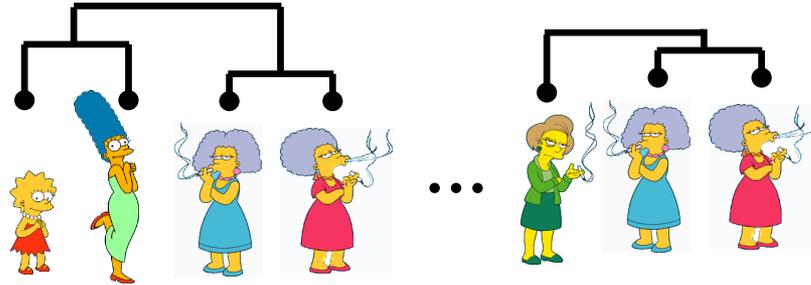
0	8	8	7	7
	0	2	4	4
		0	3	3
			0	1
				0

Bottom-Up (aglomerativo):

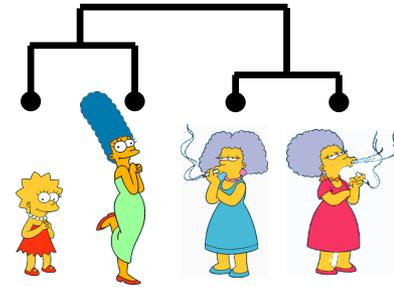
Iniciando com cada objeto em seu próprio cluster, encontra o melhor par de *clusters* para unir em um novo *cluster*. Repete até que todos os *clusters* sejam um único *cluster*.



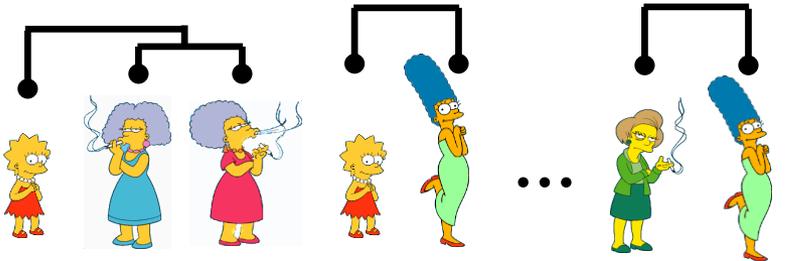
Considerar todas as uniões possíveis ...



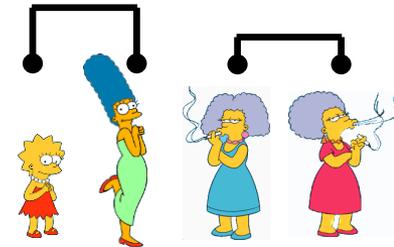
Escolher a melhor



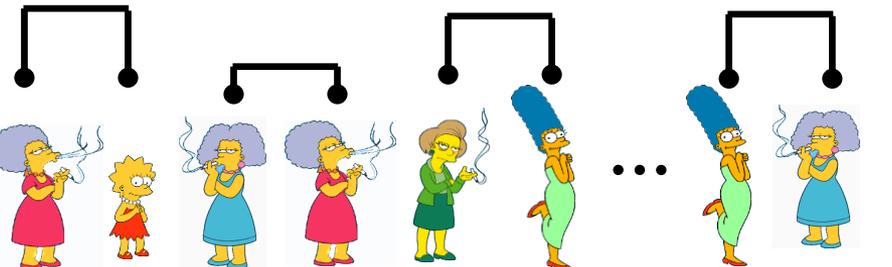
Considerar todas as uniões possíveis ...



Escolher a melhor



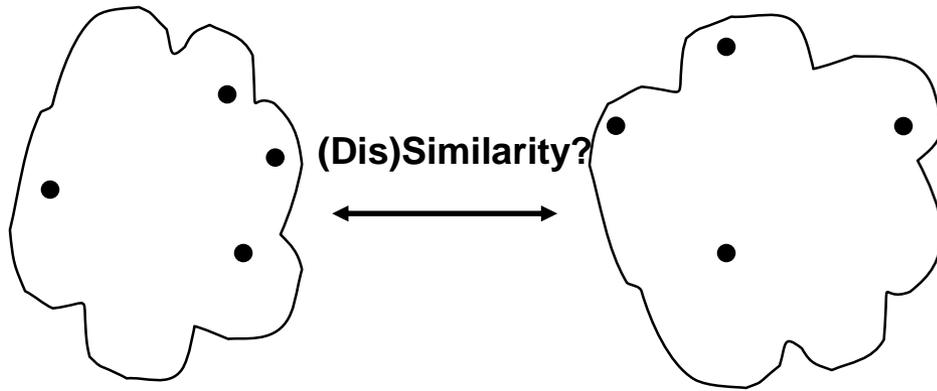
Considerar todas as uniões possíveis ...



Escolher a melhor



How to Define Inter-Cluster (Dis)Similarity

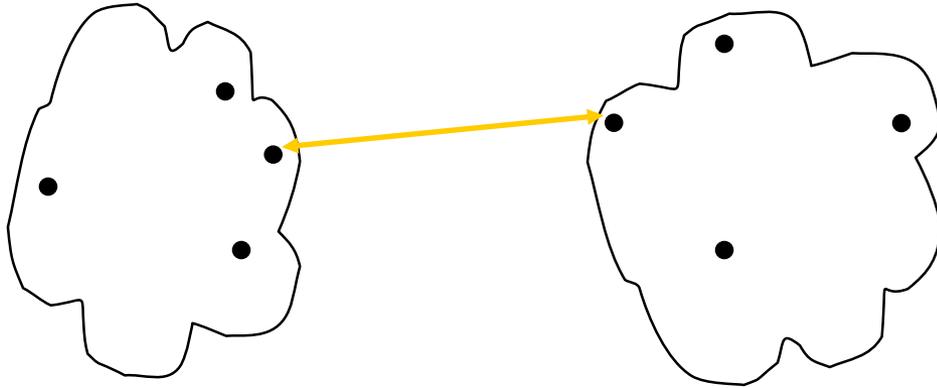


- | MIN
- | MAX
- | Group Average
- | Distance Between Centroids
- | Other methods
 - Ward's
 - ...

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

· **Proximity Matrix**

How to Define Inter-Cluster (Dis)Similarity

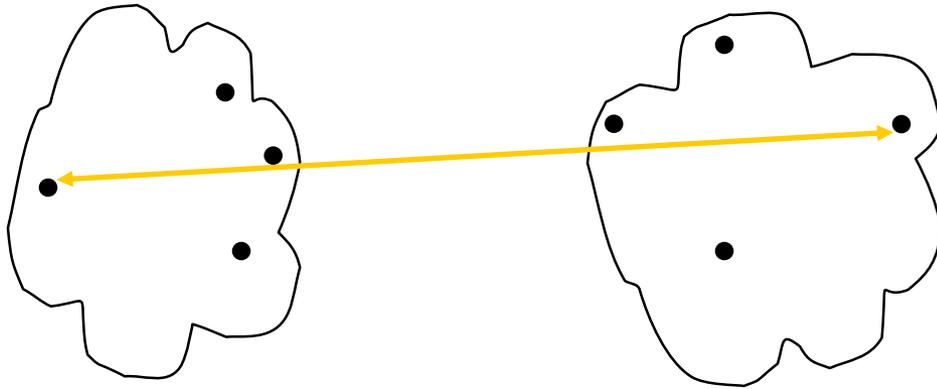


- | MIN
- | MAX
- | Group Average
- | Distance Between Centroids
- | Other methods
 - Ward's
 - ...

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

· Proximity Matrix

How to Define Inter-Cluster (Dis)Similarity

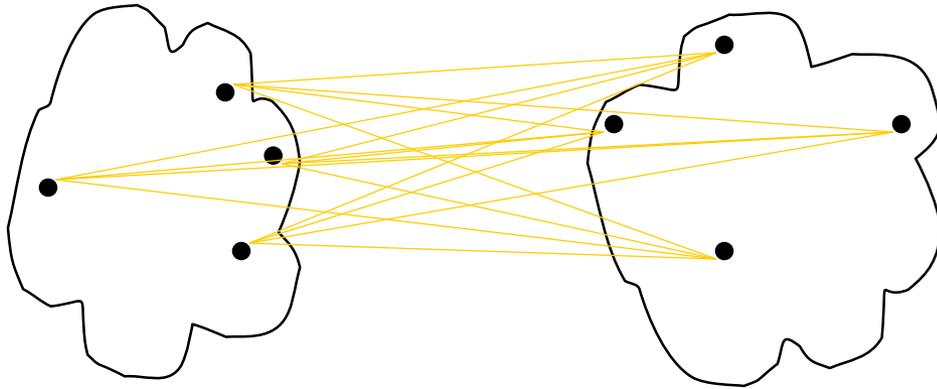


- | MIN
- | **MAX**
- | Group Average
- | Distance Between Centroids
- | Other methods
 - Ward's
 - ...

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						
.						
.						

· **Proximity Matrix**

How to Define Inter-Cluster (Dis)Similarity

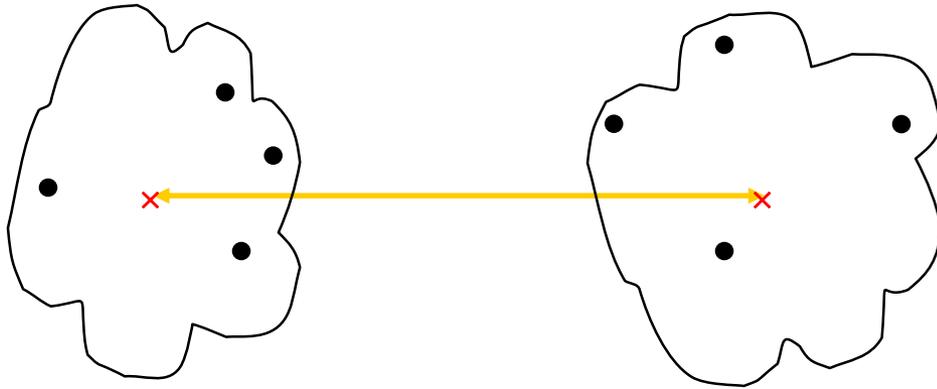


- | MIN
- | MAX
- | **Group Average**
- | Distance Between Centroids
- | Other methods
 - Ward's
 - ...

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

· Proximity Matrix

How to Define Inter-Cluster (Dis)Similarity



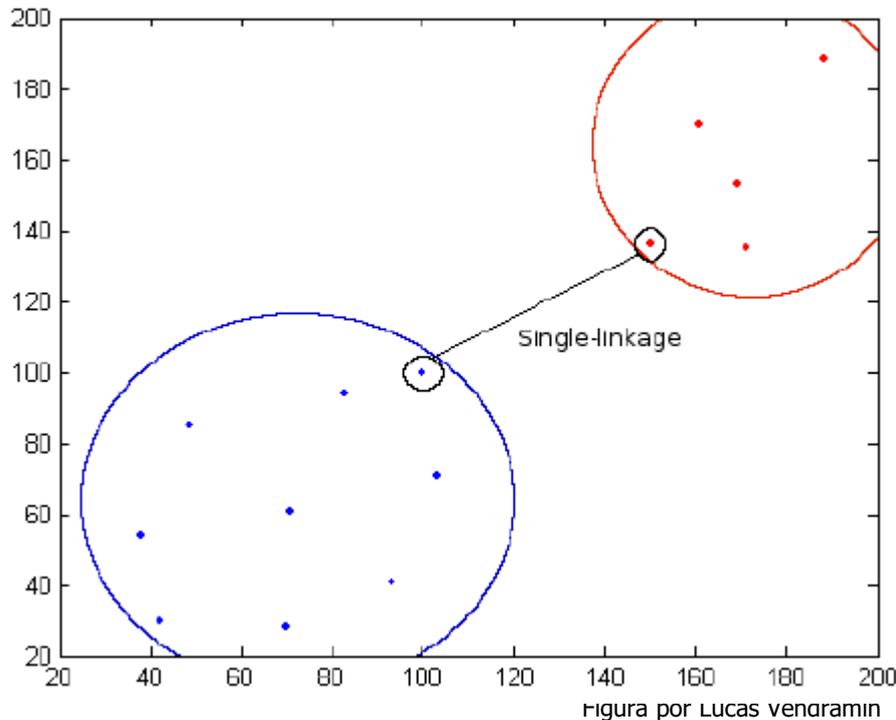
- | MIN
- | MAX
- | Group Average
- | **Distance Between Centroids**
- | Other methods
 - Ward's
 - ...

	p1	p2	p3	p4	p5	...
p1						
p2						
p3						
p4						
p5						

· **Proximity Matrix**

Como comparar os clusters?

- ***Single Linkage***, Min, ou Vizinho mais Próximo :
 - Dissimilaridade entre *clusters* é dada pela menor dissimilaridade entre 2 objetos (um de cada *cluster*)



single link (Florek, 1951; Sneath, 1957)

Originalmente baseado em **Grafos**:
menor aresta entre dois vértices
de subconjuntos distintos

Propriedade Útil

- Propriedade da Função Mínimo (min):
 - $\min\{\mathbf{D}\} = \min\{ \min\{\mathbf{D}_1\} , \min\{\mathbf{D}_2\} \}$
 - \mathbf{D} , \mathbf{D}_1 e \mathbf{D}_2 são conjuntos de valores reais tais que $\mathbf{D}_1 \cup \mathbf{D}_2 = \mathbf{D}$
 - Exemplo:
 - $\min\{10, -3, 0, 100\} = \min \{ \min\{10, -3\}, \min\{0, 100\} \} = -3$
 - Propriedade vale recursivamente (para $\min\{\mathbf{D}_1\}$ e $\min\{\mathbf{D}_2\}$)
- Utilidade para *Single-Linkage*
 - Dada a distância entre os grupos \mathbf{A} e \mathbf{B} e entre \mathbf{A} e \mathbf{C}
 - É trivial calcular a distância entre \mathbf{A} e $(\mathbf{B} \cup \mathbf{C})$.

Exemplo de Single Linkage: Método de Johnson (1967)

- Consideremos a seguinte matriz de distâncias iniciais (\mathbf{D}_1) entre 5 objetos $\{1,2,3,4,5\}$. Qual par de objetos será escolhido para formar o 1º *cluster*?

$$\mathbf{D}_1 = \begin{matrix} & \begin{matrix} 1 & 2 & 3 & 4 & 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & & & & \\ & 2 & & & \\ & 6 & 5 & & \\ & 10 & 9 & 4 & \\ & 9 & 8 & 5 & 3 & 0 \end{bmatrix} \end{matrix}$$

- A menor distância entre objetos é $d_{12}=d_{21}=2$, indicando que estes dois objetos serão unidos em um *cluster*. Na seqüência, calcula-se:

$$d_{(12)3} = \min\{d_{13}, d_{23}\} = d_{23} = 5;$$

$$d_{(12)4} = \min\{d_{14}, d_{24}\} = d_{24} = 9;$$

$$d_{(12)5} = \min\{d_{15}, d_{25}\} = d_{25} = 8;$$

- Desta forma, obtém-se uma nova matriz de distâncias (\mathbf{D}_2), que será usada na próxima etapa do agrupamento hierárquico:

$$\mathbf{D}_2 = \begin{matrix} & 12 & \begin{bmatrix} 0 \\ 3 & 5 & 0 \\ 4 & 9 & 4 & 0 \\ 5 & 8 & 5 & \boxed{3} & 0 \end{bmatrix} \end{matrix}$$

- Qual o novo *cluster* a ser formado?
- Unindo os objetos **4** e **5** obtemos três clusters: {1,2}, {4,5}, {3}
- Como $d_{(12)3}$ já está calculada, calculamos na sequência:

$$d_{(12)(45)} = \min\{d_{(12)(4)}, d_{(12)(5)}\} = d_{(12)(5)} = 8$$

$$d_{(45)3} = \min\{d_{43}, d_{53}\} = d_{43} = 4$$

obtendo a seguinte matriz:

$$\mathbf{D}_3 = \begin{matrix} & 12 & \begin{bmatrix} 0 \\ 3 & 5 & 0 \\ 45 & 8 & \boxed{4} & 0 \end{bmatrix} \end{matrix}$$

* Unir *cluster* {3} com {4,5};

* Finalmente, unir todos os *clusters* em um único *cluster*

- A sequência de partições obtidas neste exemplo é, portanto:

$$\{ (1), (2), (3), (4), (5) \} \rightarrow \{ (1, 2), (3), (4), (5) \} \rightarrow$$

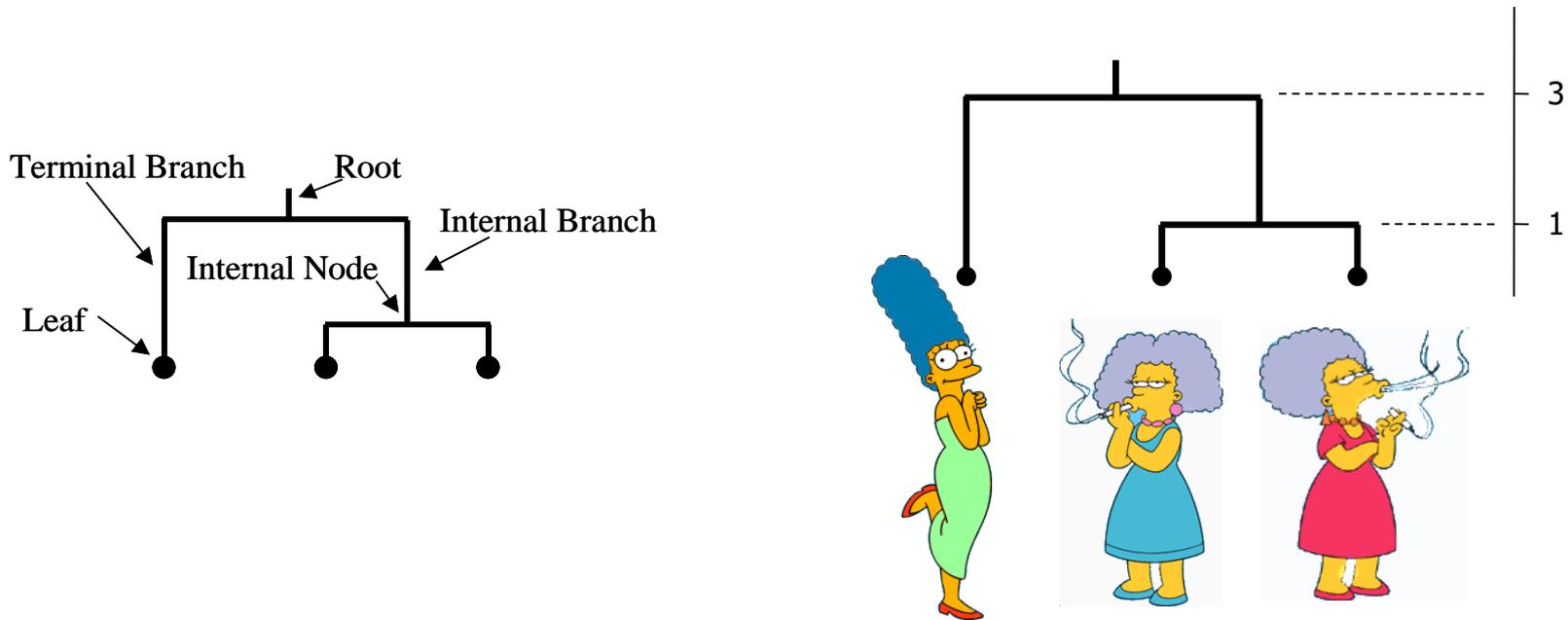
$$\{ (1, 2), (3), (4, 5) \} \rightarrow \{ (1, 2), (3, 4, 5) \} \rightarrow \{ (1, 2, 3, 4, 5) \}$$

- **Homework:** Para *single link*, a dissimilaridade entre 2 *clusters* pode ser computada a partir da matriz atualizada na iteração anterior, sem necessidade da matriz original

- Isso vale devido à propriedade da função *min* (slide oculto)
- No nosso exemplo, simplificamos o cálculo de $d_{(12)(45)}$ como $\min\{d_{(12)(4)}, d_{(12)(5)}\}$ fazendo uso daquela propriedade:
 - $\min\{d_{(12)(4)}, d_{(12)(5)}\} = \min\{9, 8\} = \min\{d_{14}, d_{24}, d_{15}, d_{25}\}$

Dendrograma

Dendrograma: Hierarquia + Dissimilaridades entre Clusters

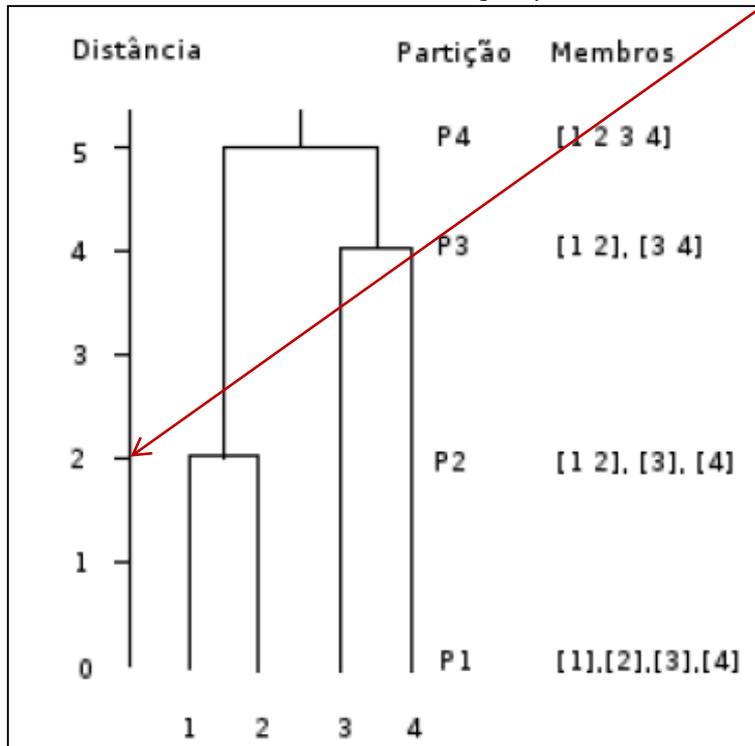


* A dissimilaridade entre dois clusters (possivelmente **singletons**) é representada como a altura do nó interno mais baixo compartilhado

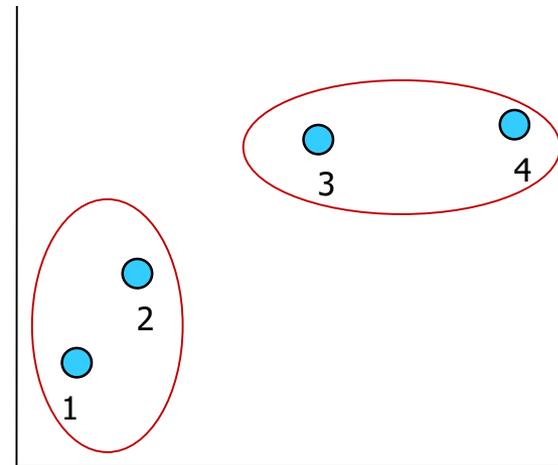
Exemplo de Dendrograma

$$\mathbf{D} = \begin{bmatrix} 1 & 0 & 2 & 7 & 13 \\ 2 & 2 & 0 & 5 & 10 \\ 3 & 7 & 5 & 0 & 4 \\ 4 & 13 & 10 & 4 & 0 \end{bmatrix}$$

Figura por Lucas Vendramin

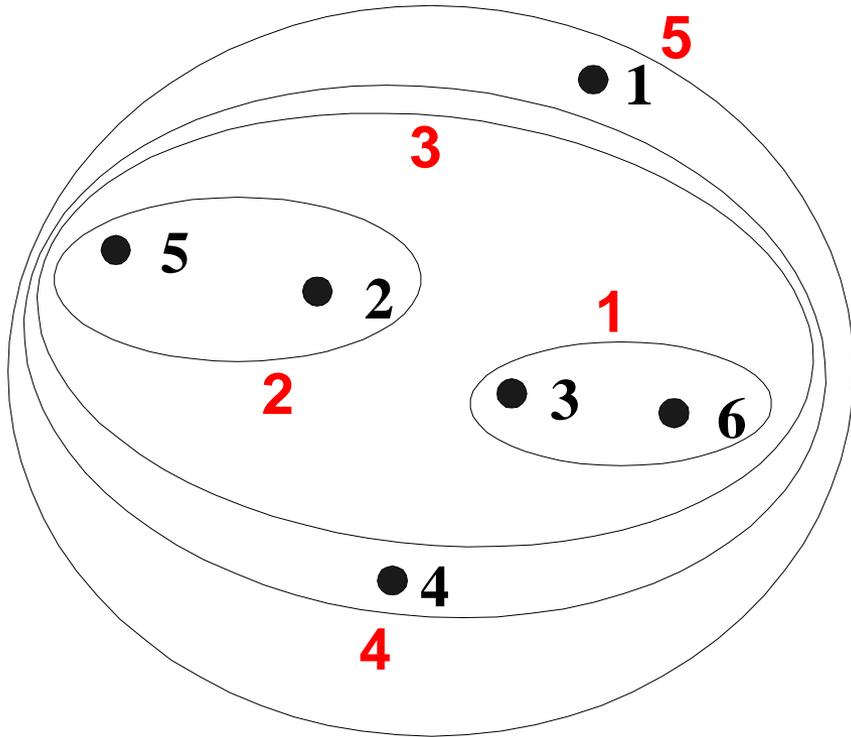


Dendrograma

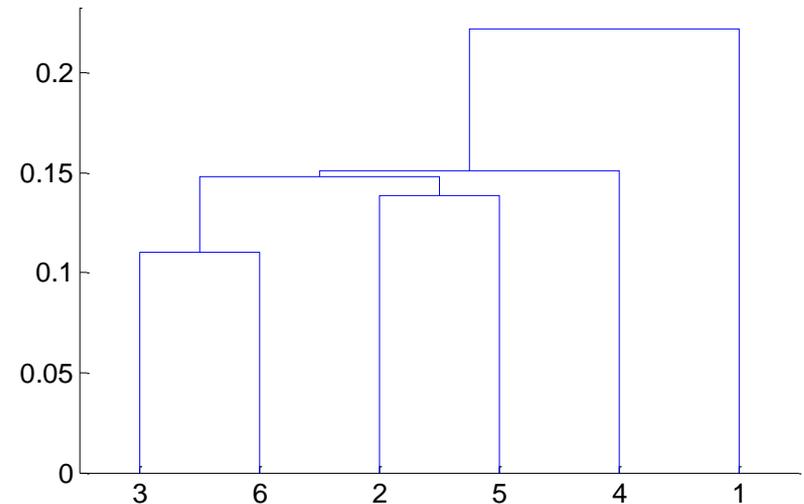


uma das partições aninhadas

Outro Exemplo de Dendrograma



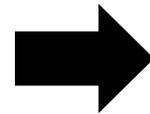
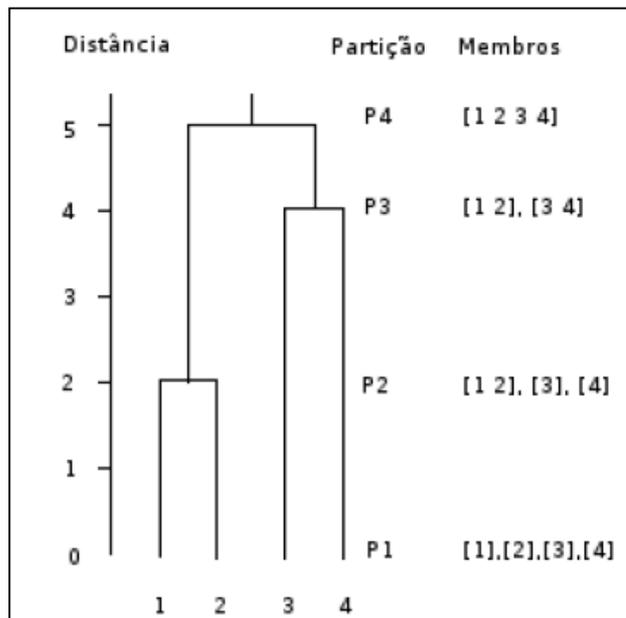
Nested Clusters



Dendrogram

Cophenetic Matrix

- Matriz com as dissimilaridades que levaram à união de cada par de objetos na base de dados. Exemplo:



$$C_p = \begin{bmatrix} 0 & 2 & 5 & 5 \\ 2 & 0 & 5 & 5 \\ 5 & 5 & 0 & 4 \\ 5 & 5 & 4 & 0 \end{bmatrix}$$

- Esta matriz é importante para a validação de agrupamentos hierárquicos (tópico a ser discutido posteriormente no curso)

Exercício:

- Obtenha o dendrograma completo para o exemplo visto de execução do *single linkage* (matriz de distâncias abaixo)

$$\mathbf{D}_1 = \begin{matrix} & 1 & 2 & 3 & 4 & 5 \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & & & & \\ 2 & 0 & & & \\ 6 & 5 & 0 & & \\ 10 & 9 & 4 & 0 & \\ 9 & 8 & 5 & 3 & 0 \end{bmatrix} \end{matrix}$$

- Apresente também a *cophenetic matrix* correspondente

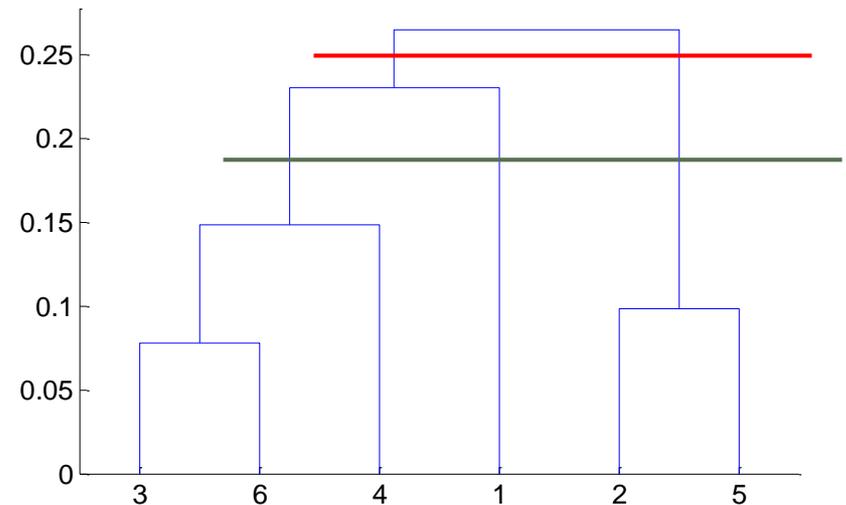
Dendrogramas e Partições

- Partições são obtidas via **cortes** no dendrograma
 - cortes horizontais
 - no. de grupos da partição = no. de interseções

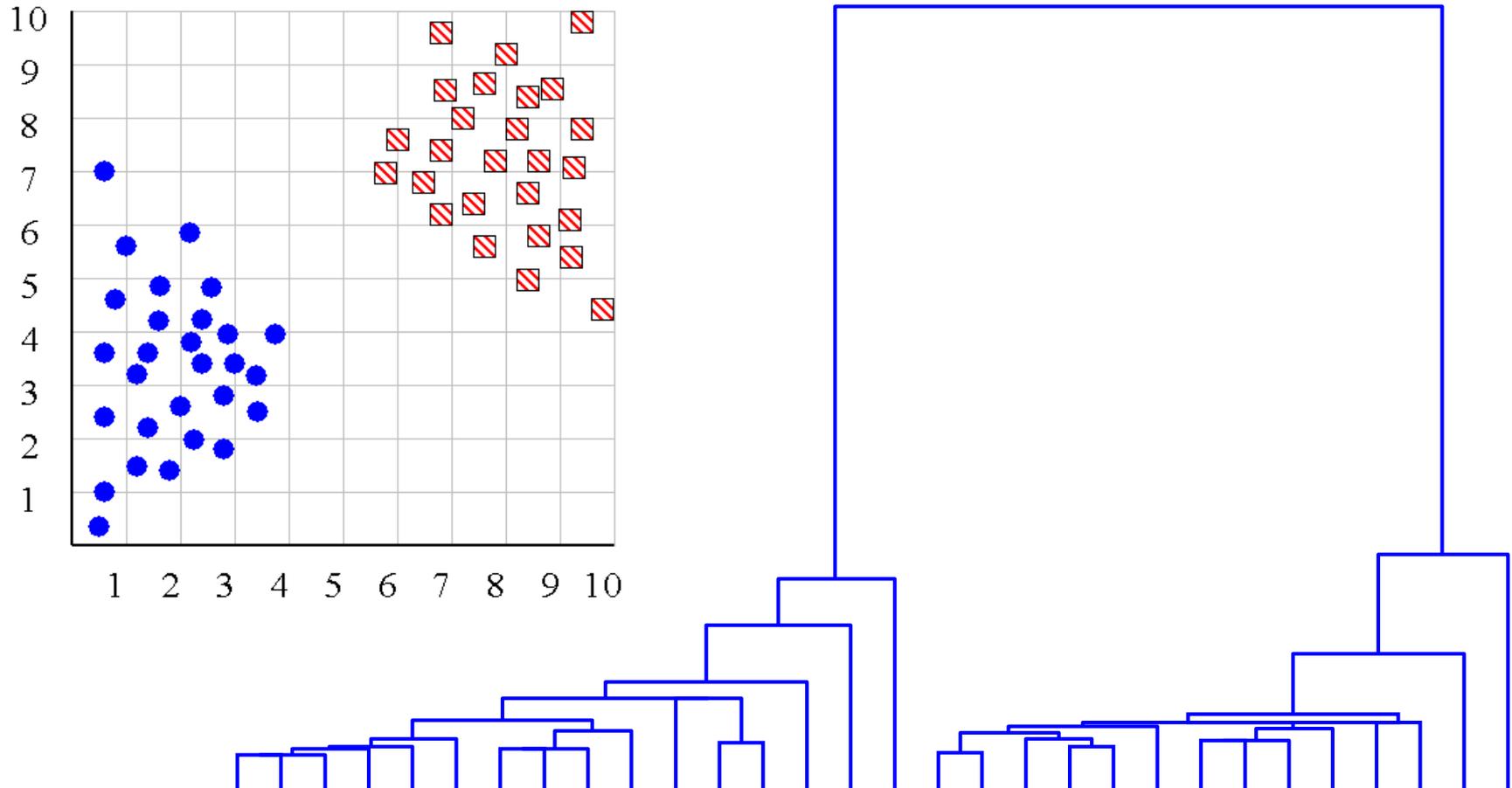
➤ Exemplos:

$$P_2 = \{ (\mathbf{x}_1, \mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$$

$$P_1 = \{ (\mathbf{x}_1), (\mathbf{x}_3, \mathbf{x}_4, \mathbf{x}_6), (\mathbf{x}_2, \mathbf{x}_5) \}$$



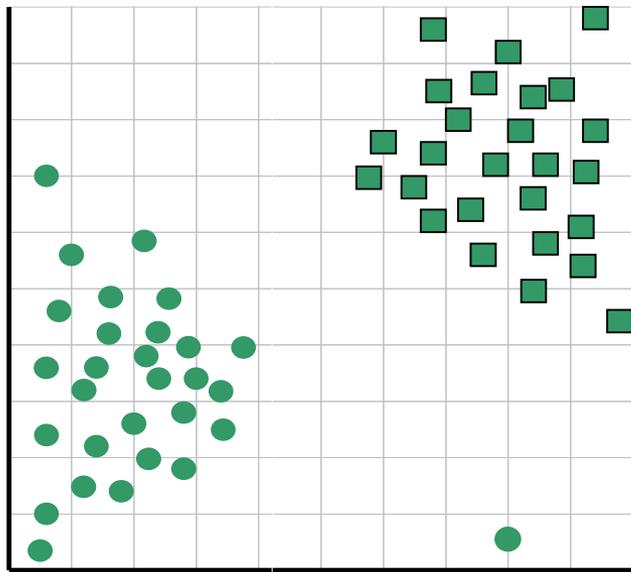
Pode-se examinar o dendrograma para tentar estimar o número *natural* de *clusters*. No caso abaixo, existem duas sub-árvores bem separadas, sugerindo dois *clusters*.



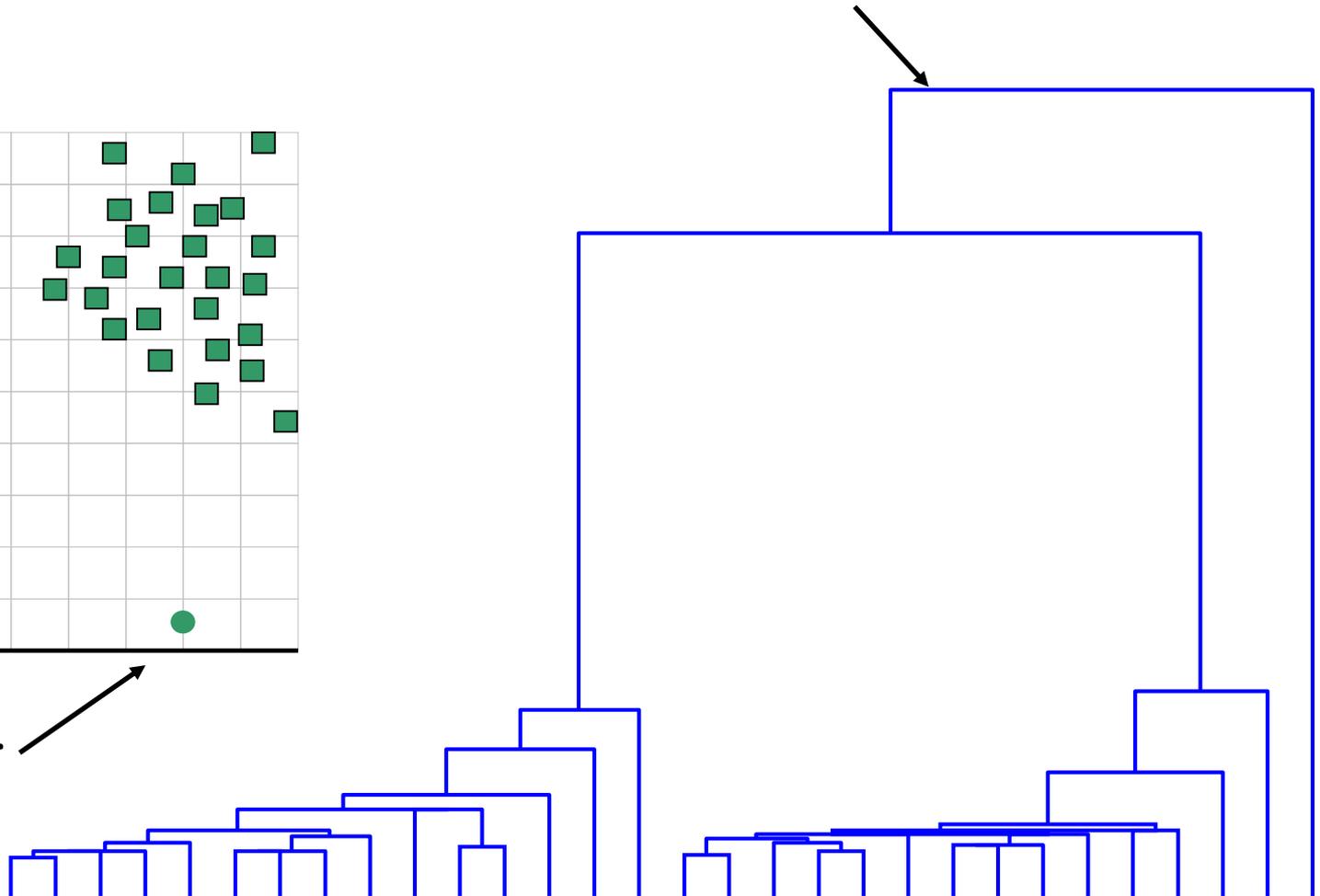
➤ Infelizmente, na prática, as distinções não são tão simples...

Pode-se usar o dendrograma para tentar detectar *outliers*:

Ramo isolado sugere que o objeto é muito diferente dos demais.



Outlier



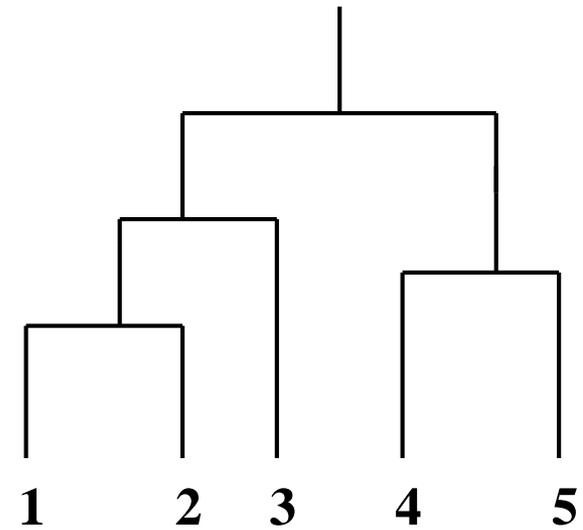
Como Plotar o Dendrograma ?

- Para muitos objetos, ilustração manual é inviável
- Gerar o gráfico automaticamente demanda ordenar de forma apropriada os objetos no eixo horizontal
- Algoritmo Recursivo Simples:
 - Iniciar com o topo da hierarquia (grupo único)
 - Dividir o eixo horizontal em 2 subintervalos e colocar em cada um os objetos de cada um dos 2 grupos que derivam do grupo único
 - Executar recursivamente o passo anterior para cada subintervalo

Voltando ao Single Linkage (Min)...

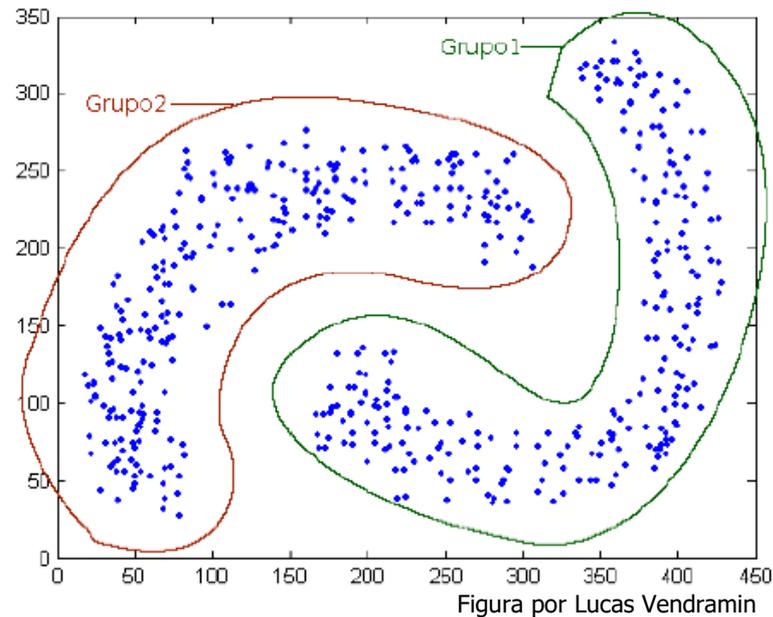
- | Similarity of two clusters is based on the two most similar (closest) points in the clusters
 - Determined by **one pair of points**
 - i.e., by **one link** in the **proximity graph**

	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



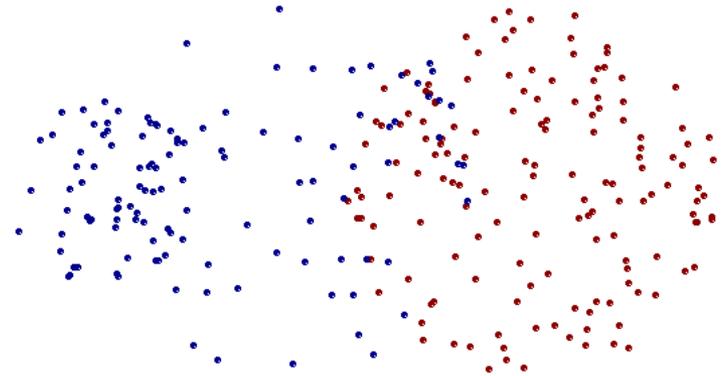
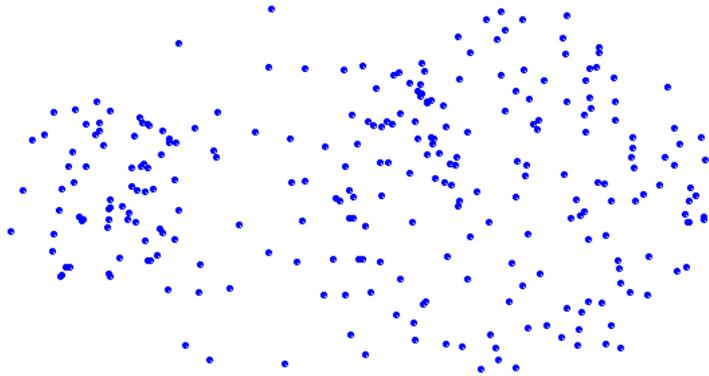
Strength of MIN (Single Link)

- Can handle non-elliptical shapes



Main Limitations of MIN

- Sensitive to noise and outliers



- ***Complete Linkage***, Max, ou Vizinho mais Distante:

- Dissimilaridade entre *clusters* é dada pela maior dissimilaridade entre dois objetos (um de cada *cluster*)

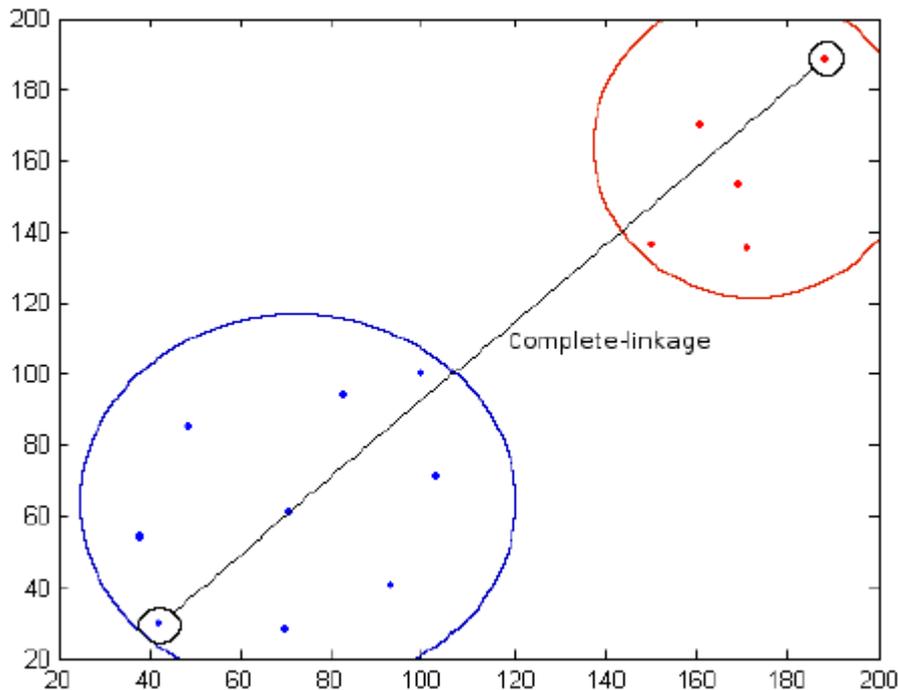


Figura por Lucas Vendramin

complete link (Sorensen, 1948)

Originalmente baseado em **Grafos**:
maior aresta entre dois vértices de subconjuntos distintos

Propriedade Útil

- Propriedade da Função Máximo (max):
 - $\max\{\mathbf{D}\} = \max\{ \max\{\mathbf{D}_1\} , \max\{\mathbf{D}_2\} \}$
 - \mathbf{D} , \mathbf{D}_1 e \mathbf{D}_2 são conjuntos de valores reais tais que $\mathbf{D}_1 \cup \mathbf{D}_2 = \mathbf{D}$
 - Exemplo:
 - $\max\{10, -3, 0, 100\} = \max\{ \max\{10, -3\}, \max\{0, 100\} \} = 100$
 - Propriedade vale recursivamente (para $\max\{\mathbf{D}_1\}$ e $\max\{\mathbf{D}_2\}$)
- Utilidade para *Complete-Linkage*
 - Dada a distância entre os grupos \mathbf{A} e \mathbf{B} e entre \mathbf{A} e \mathbf{C}
 - É trivial calcular a distância entre \mathbf{A} e $(\mathbf{B} \cup \mathbf{C})$.

- Seja a seguinte matriz de distâncias iniciais (\mathbf{D}_1) entre 5 objetos :

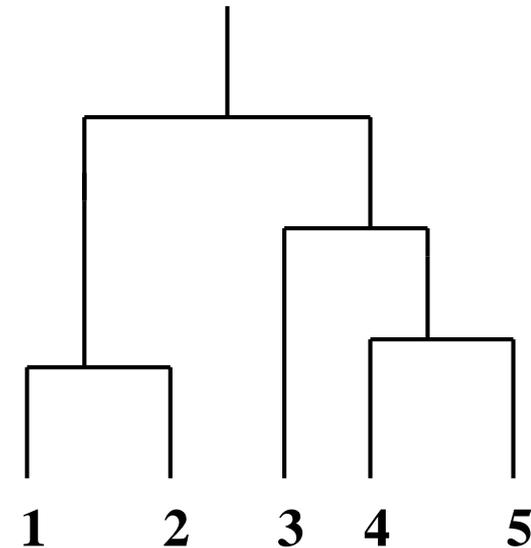
$$\mathbf{D}_1 = \begin{matrix} & \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} \\ \begin{matrix} 1 \\ 2 \\ 3 \\ 4 \\ 5 \end{matrix} & \begin{bmatrix} 0 & & & & \\ & 2 & 0 & & \\ & 6 & 5 & 0 & \\ & 10 & 9 & 4 & 0 \\ & 9 & 8 & 5 & 3 & 0 \end{bmatrix} \end{matrix}$$

- Exercício: executar o *complete linkage* através de sucessivas atualizações da matriz de distâncias (método de Johnson).

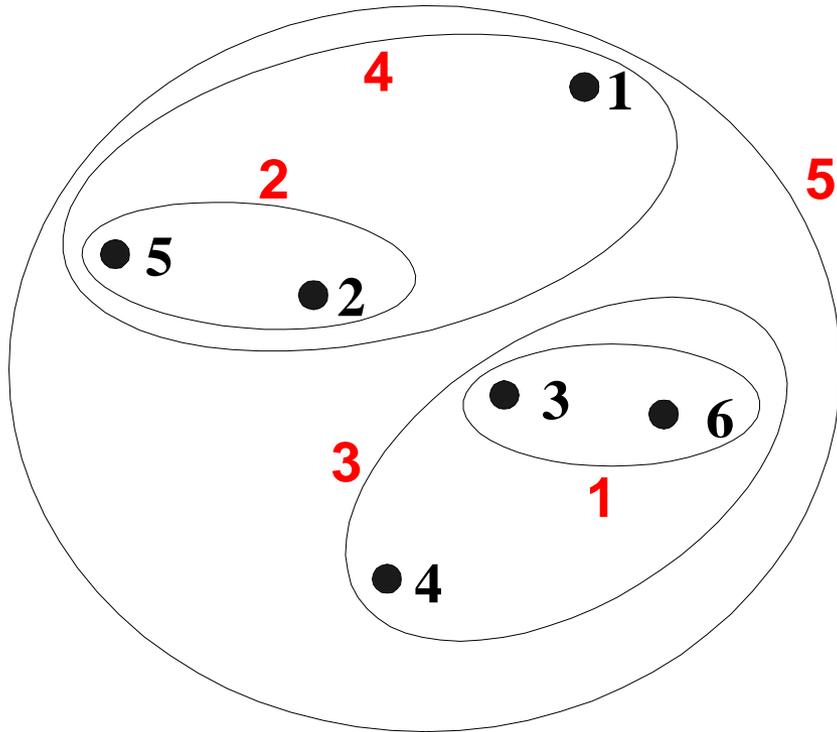
Cluster Similarity: MAX or Complete Linkage

- | Similarity of two clusters is based on the two least similar (most distant) points in the clusters
 - Determined by **one pair of points**

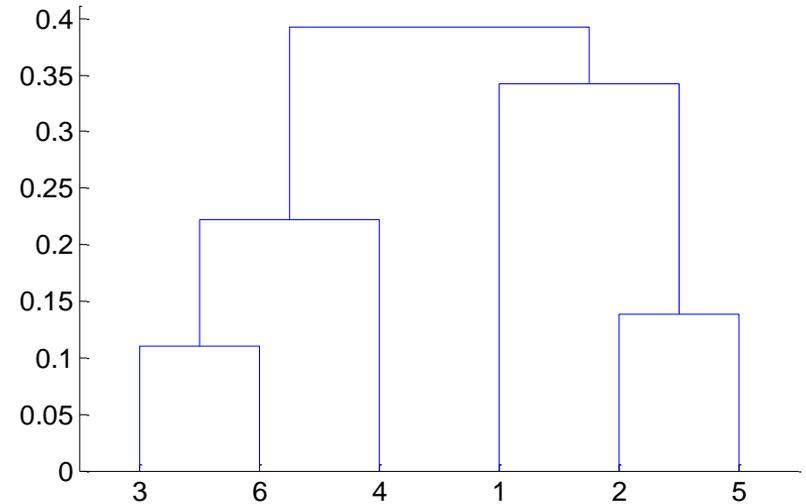
	I1	I2	I3	I4	I5
I1	1.00	0.90	0.10	0.65	0.20
I2	0.90	1.00	0.70	0.60	0.50
I3	0.10	0.70	1.00	0.40	0.30
I4	0.65	0.60	0.40	1.00	0.80
I5	0.20	0.50	0.30	0.80	1.00



Hierarchical Clustering: MAX

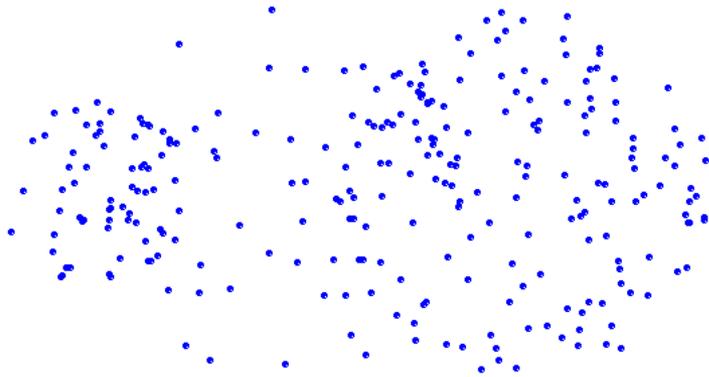


Nested Clusters

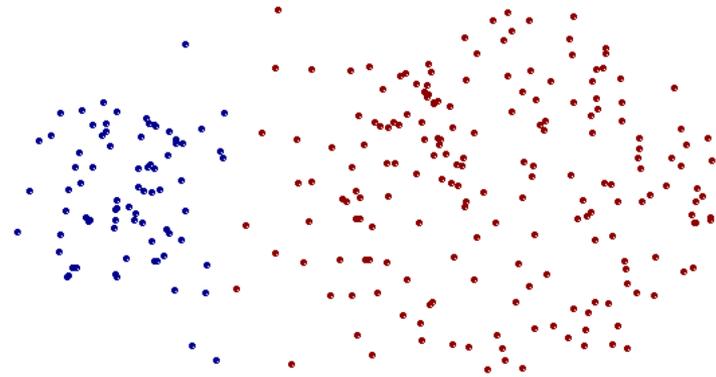


Dendrogram

Strength of MAX



Original Points

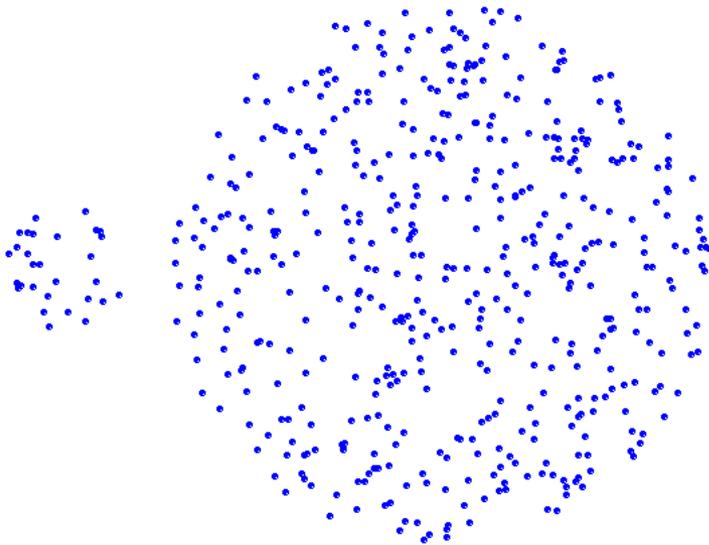


Two Clusters

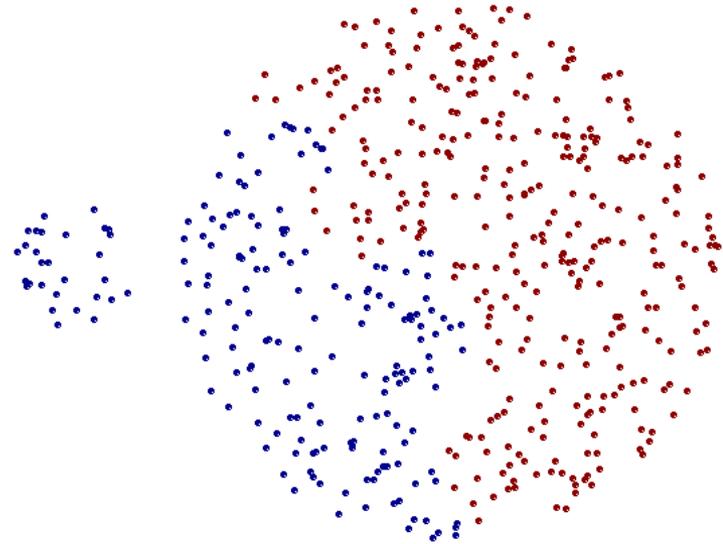
- **Less susceptible to noise and outliers**

Main Limitations of MAX

Original Points



Two Clusters



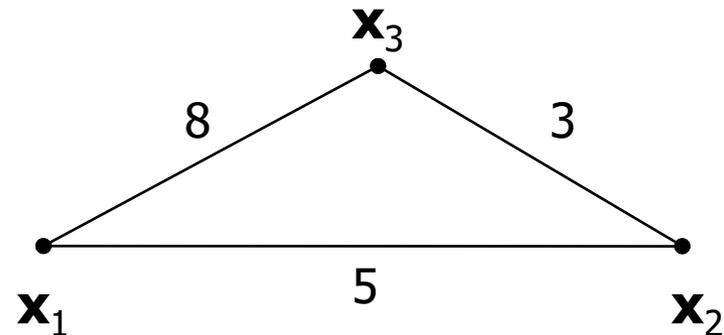
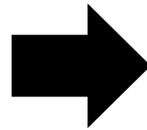
- **Tends to break large clusters**
- **Biased towards globular clusters**

Vinculação Simples/Completa sob Perspectiva de Teoria dos Grafos

Grafo de Proximidades

- Grafo (**ponderado**, sem laços e sem múltiplas arestas) no qual:
 - vértices representam os objetos da base de dados
 - arestas representam as (dis)similaridades entre pares de objetos
- Exemplo (3 objetos):

$$\mathbf{D} = \begin{matrix} 1 \\ 2 \\ 3 \end{matrix} \begin{bmatrix} 0 & 5 & 8 \\ 5 & 0 & 3 \\ 8 & 3 & 0 \end{bmatrix}$$



➤ Por simplicidade, assumamos uma matriz de distâncias em escala ordinal (sem empates)

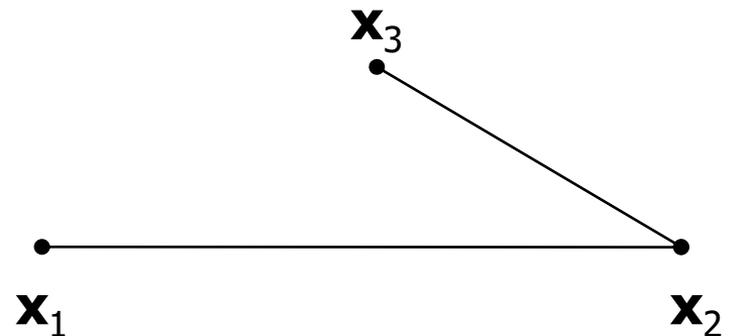
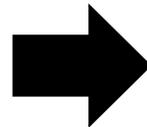
➤ Exemplo:

$$\mathbf{D} = \begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{matrix} \begin{bmatrix} 0 & 6 & 8 & 2 & 7 \\ 6 & 0 & 1 & 5 & 3 \\ 8 & 1 & 0 & 10 & 9 \\ 2 & 5 & 10 & 0 & 4 \\ 7 & 3 & 9 & 4 & 0 \end{bmatrix}$$

Grafo de Limiar $G(\nu)$

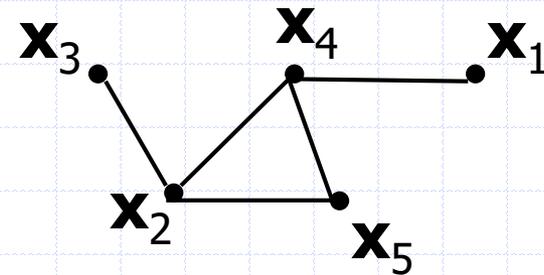
- Grafo (**não ponderado**, sem laços e sem múltiplas arestas) no qual:
 - vértices de $G(\nu)$ representam os objetos da base de dados
 - arestas (links) de $G(\nu)$ conectam pares de objetos que possuem dissimilaridade menor ou igual a um limiar ν
 - ou similaridade maior ou igual a um limiar ν
- Exemplo Simples: $G(5)$

$$\mathbf{D} = \begin{matrix} 1 & \begin{bmatrix} 0 & 5 & 8 \end{bmatrix} \\ 2 & \begin{bmatrix} 5 & 0 & 3 \end{bmatrix} \\ 3 & \begin{bmatrix} 8 & 3 & 0 \end{bmatrix} \end{matrix}$$



- $G(v)$ define uma **relação binária** para qualquer n° real v
 - subconjunto do Produto Cartesiano $\mathbf{X} \times \mathbf{X}$, onde $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- Exemplo para $v = 5$:

$$\mathbf{D} = \begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{matrix} \begin{bmatrix} 0 & 6 & 8 & 2 & 7 \\ 6 & 0 & 1 & 5 & 3 \\ 8 & 1 & 0 & 10 & 9 \\ 2 & 5 & 10 & 0 & 4 \\ 7 & 3 & 9 & 4 & 0 \end{bmatrix}$$



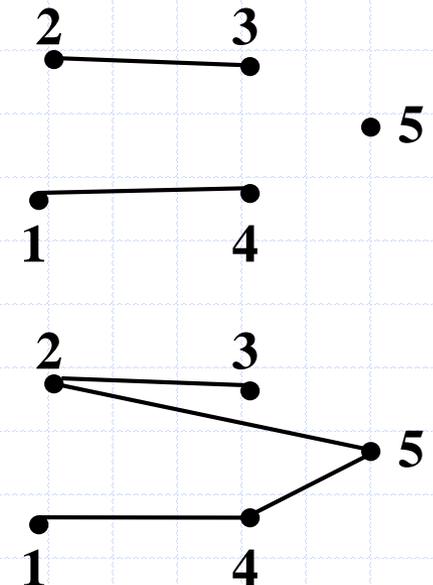
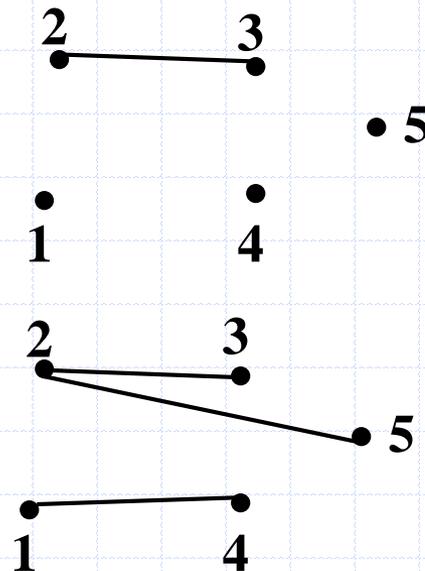
$$\mathbf{R} = \begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{matrix} \begin{bmatrix} 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 & 1 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Algoritmo Vinculação Simples (*single link*):

- 1) Iniciar com $G(0)$: cada objeto formando um grupo. $k \leftarrow 1$
- 2) Formar o grafo de limiar $G(k)$. Se o n° de componentes conexos em $G(k)$ for menor do que o n° de grupos corrente, re-nomear cada um dos componentes como sendo um grupo
- 3) Se $G(k)$ formar um único componente conexo, parar. Senão, fazer $k \leftarrow k + 1$ e voltar ao passo 2

Exemplo:

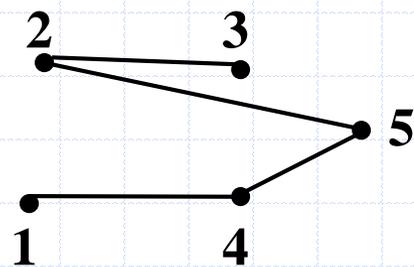
$$\mathbf{D} = \begin{matrix} \mathbf{x}_1 \\ \mathbf{x}_2 \\ \mathbf{x}_3 \\ \mathbf{x}_4 \\ \mathbf{x}_5 \end{matrix} \begin{bmatrix} 0 & 6 & 8 & \textcircled{2} & 7 \\ 6 & 0 & \textcircled{1} & 5 & \textcircled{3} \\ 8 & 1 & 0 & 10 & 9 \\ 2 & 5 & 10 & 0 & \textcircled{4} \\ 7 & 3 & 9 & 4 & 0 \end{bmatrix}$$



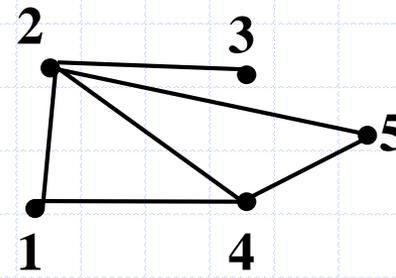
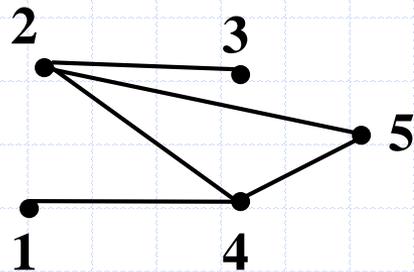
Algoritmo Vinculação Completa:

- 1) Iniciar com $G(0)$: cada objeto formando um grupo. $k \leftarrow 1$
- 2) Formar o grafo de limiar $G(k)$. Se 2 dos grupos correntes formam um *clique* (subgrafo completo) em $G(k)$, unir tais grupos
- 3) Se $k = N(N-1)/2$, o que implica que $G(k)$ é um grafo completo, parar. Caso contrário, $k \leftarrow k + 1$ e voltar ao passo 2

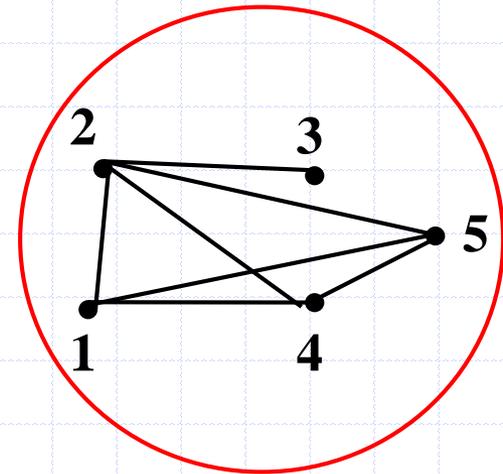
Exemplo (continuação):



grupos correntes: (2,3) (1,4) (5)



clique de grupos correntes $\rightarrow (1, 4, 5)$



* Dois grupos: pode-se interromper o processo...

Hierarquias Resultantes:

Vinculação **simples**:

$\{(2,3),1,4,5\}$

$\{(2,3),(1,4),5\}$

$\{(2,3,5),(1,4)\}$

Vinculação **completa**:

$\{(2,3),1,4,5\}$

$\{(2,3),(1,4),5\}$

$\{(2,3),(1,4,5)\}$

Notem que x_5 está em grupos diferentes

- Validação (a ser estudada posteriormente no curso) poderá nos auxiliar nesses casos

Observações

- **Matrizes não ordinais:**
 - ordena-se os valores de dissimilaridade e faz-se k assumir esses valores (ao invés de inteiros)
 - inicia em $k = 0$ e depois assume os valores de dissimilaridade do menor para o maior
 - os valores de k para os quais ocorre a união de dois grupos são armazenados para depois traçar o dendrograma

Observações

- **Empates** na matriz de proximidades:
 - devem ser resolvidos arbitrariamente
 - é simples verificar que **single linkage é invariante!**
 - 2 uniões candidatas empatadas serão sempre feitas (e em seguida)
 - mas **complete link** pode ser fortemente afetado pela decisão...
 - **Exercício:** mostrar isso através de um exemplo

Observações

■ Propriedades:

- O algoritmo de grafo visto anteriormente deixa evidente que *single* e *complete link* são **monótonos**
 - dissimilaridade das uniões é não decrescente
 - é crescente ao longo da hierarquia se não houver empates
 - dendrograma não possui “reversões”.
- Além disso, eles são **invariantes** a qualquer **transformação monótona** da matriz de proximidades
 - ou seja, que não altera a ordem relativa dos elementos

Observações

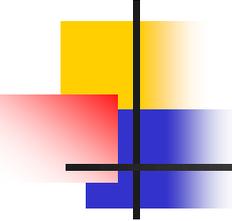
- O liberalismo de *single linkage* e o conservadorismo de *complete linkage* ficam evidentes nos grafos de limiar:
 - *single linkage* exige apenas que cada componente seja conexa e que ambas se tornem conexas após a união
 - para N objetos, é preciso apenas $N - 1$ arestas...
 - *complete linkage* exige que cada componente seja totalmente conexa (completa) e que a união também seja
 - para N objetos, demanda $N*(N - 1)/2$ arestas...
 - a aresta referente aos 2 objetos mais distantes é só a última!

Observações

- **Outras Motivações para o Estudo das Relações com Grafos:**
 - É possível modificar o algoritmo para tentar encontrar um meio termo entre liberalismo e conservadorismo
 - Por exemplo, exigindo (q, r) -conectividade:
 - todos os vértices de um grupo devem poder alcançar qualquer outro por um caminho envolvendo apenas arestas com dissimilaridade $\leq r$
 - todos os vértices de um grupo devem ser adjacentes a pelo menos q outros através de arestas com dissimilaridade $\leq r$

Observações

- **Outras Motivações para o Estudo das Relações com Grafos (cont.):**
 - À parte de questões computacionais e outras questões de cunho prático, os princípios fundamentais por trás das relações entre *clustering* e grafos são importantes
 - medidas de conectividade, p. ex., podem definir inter-relacionamentos indiretos entre objetos de um grupo, ao invés de similaridade explícita
 - Alguns algoritmos modernos de agrupamento de dados são baseados em grafos.



Referências

- Jain, A. K. and Dubes, R. C., Algorithms for Clustering Data, Prentice Hall, 1988
- Everitt, B. S., Landau, S., and Leese, M., *Cluster Analysis*, Arnold, 4th Edition, 2001.
- Tan, P.-N., Steinbach, M., and Kumar, V., *Introduction to Data Mining*, Addison-Wesley, 2006