



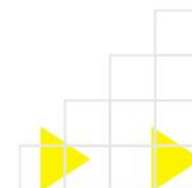
Curso 2 – CD, AM e DM

Mineração de Dados

Parte 3

Extração de Padrões
Agrupamento Hierárquico

Prof. Ricardo M. Marcacini
ricardo.marcacini@icmc.usp.br

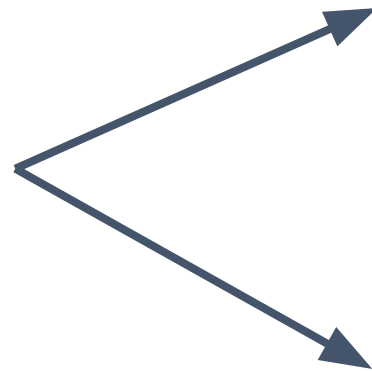
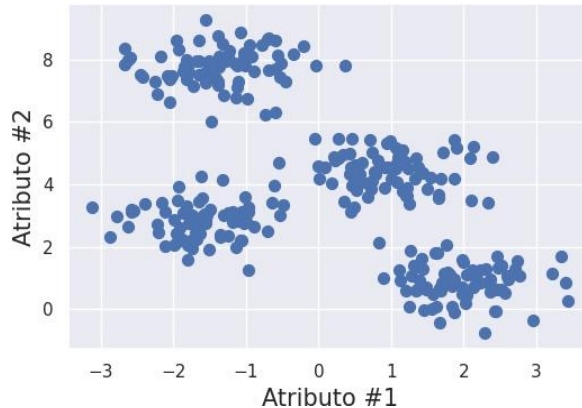


Métodos para Agrupamento de Dados

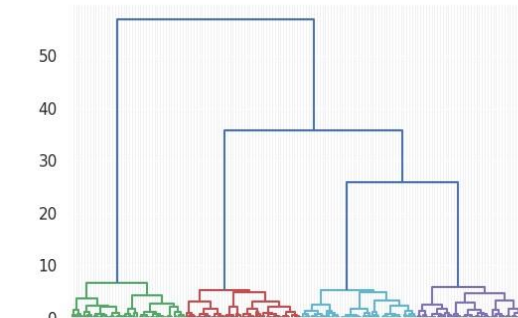


- **Particionais:** organizar dados em uma partição de k *clusters*
- **Hierárquicos:** organizar dados em uma decomposição hierárquica de *clusters* e *subclusters*

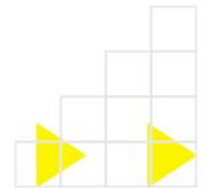
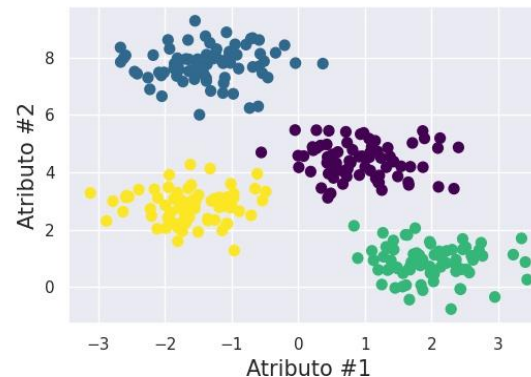
Conjunto de Dados



Agrupamento Hierárquico



Agrupamento Particional

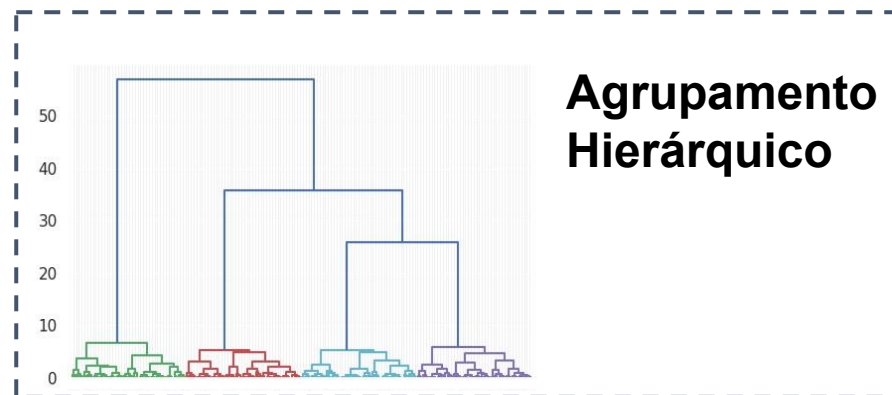
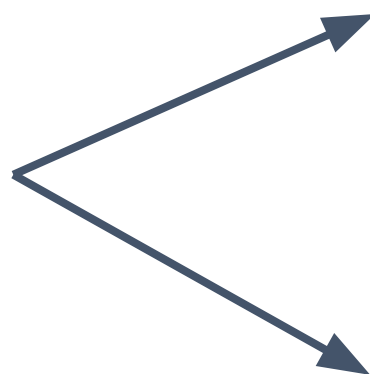
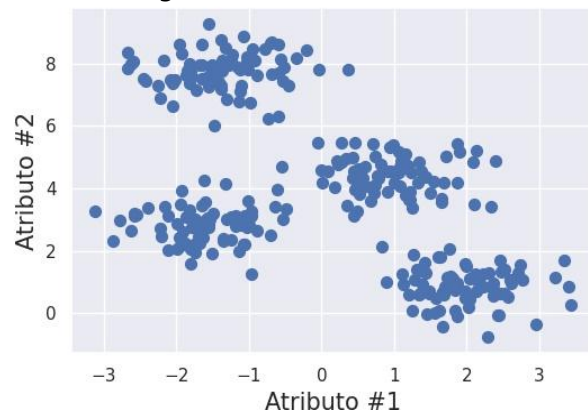


Métodos para Agrupamento de Dados

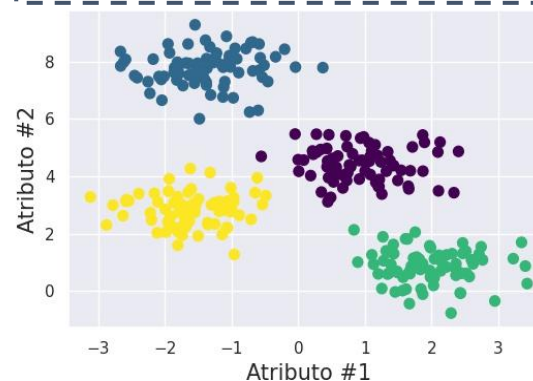


- **Particionais:** organizar dados em uma partição de k *clusters*
- **Hierárquicos:** organizar dados em uma decomposição hierárquica de *clusters* e *subclusters*

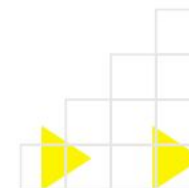
Conjunto de Dados



Agrupamento Hierárquico



Agrupamento Particional

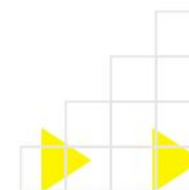
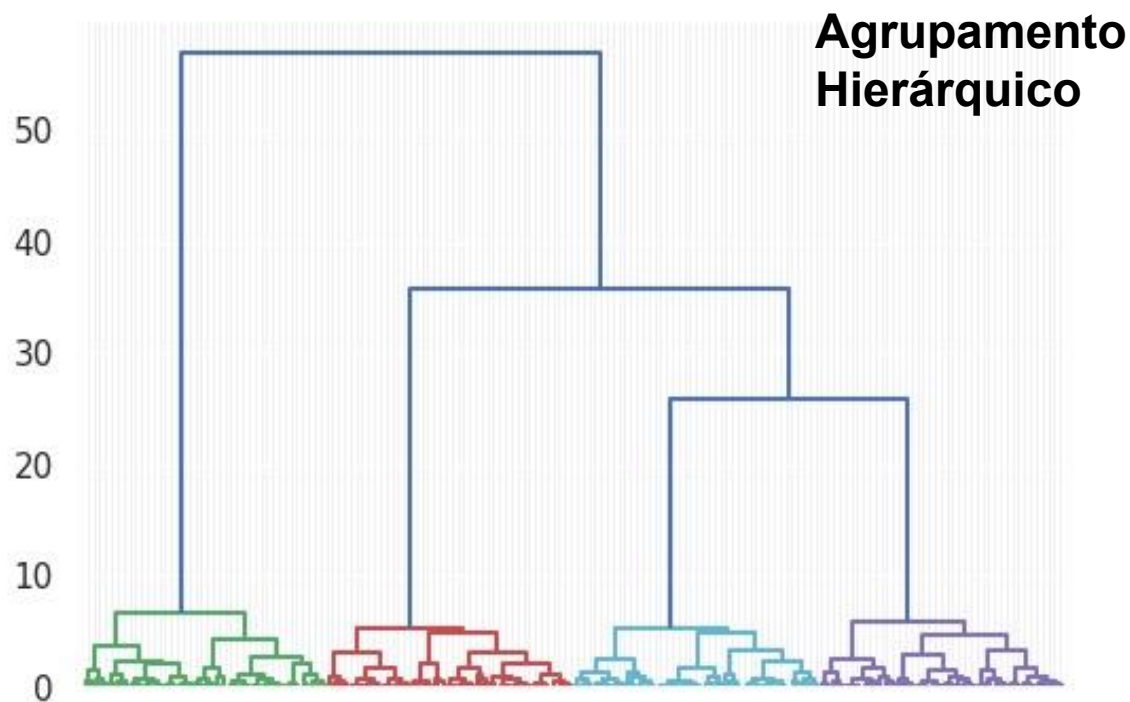
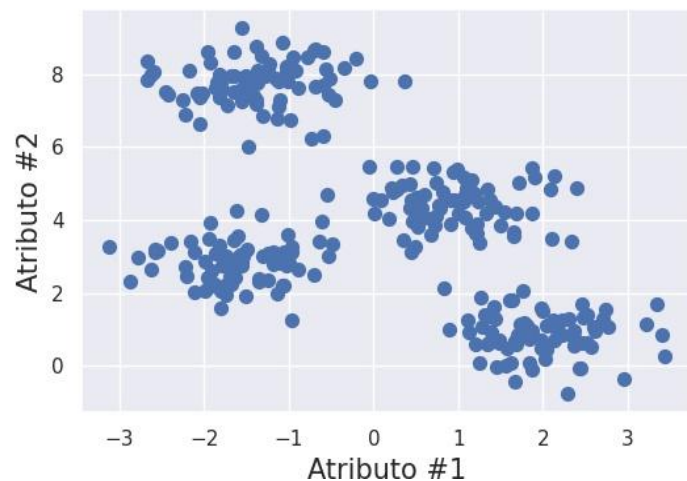


Agrupamento Hierárquico



- **Dendrograma:** diagrama com a estrutura hierárquica que representa o resultado de um agrupamento. Sumariza a formação dos *clusters* e *subclusters*.

Conjunto de Dados

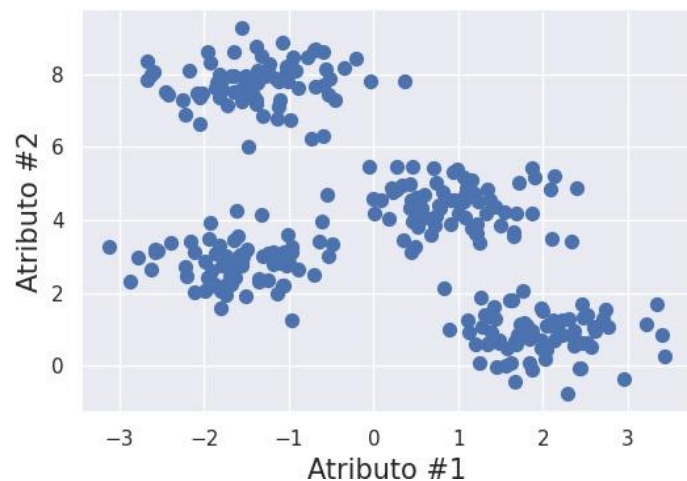


Agrupamento Hierárquico

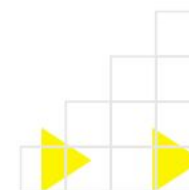
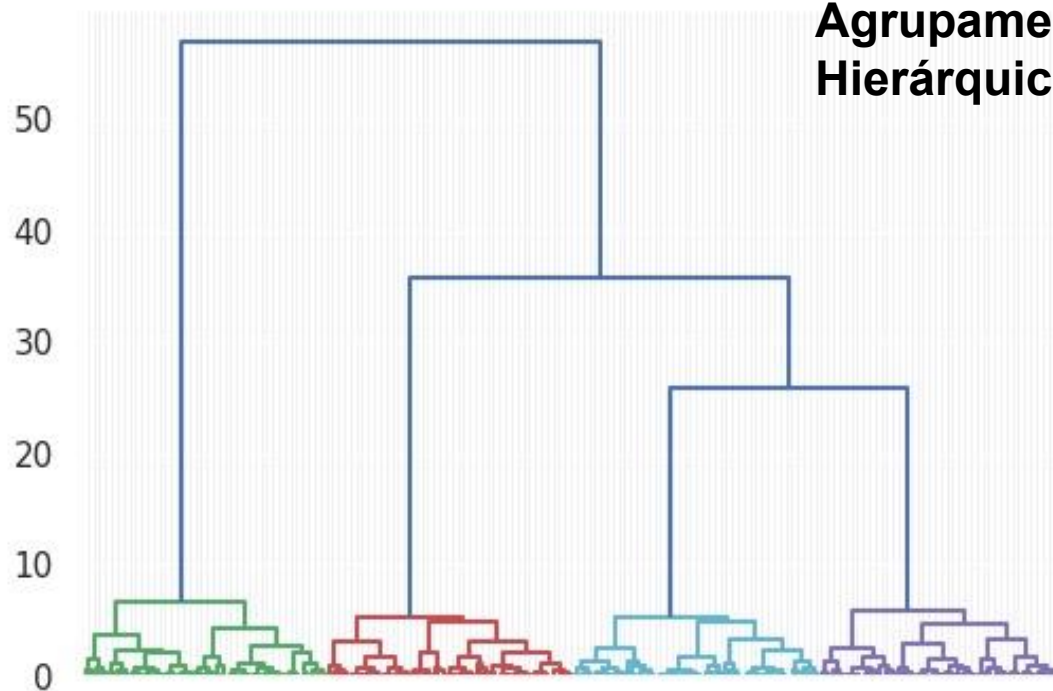


- Os objetos do conjunto de dados estão organizados no eixo x do dendrograma. A altura dos arcos indica a dissimilaridade entre objetos e grupos de objetos.

Conjunto de Dados



Agrupamento Hierárquico

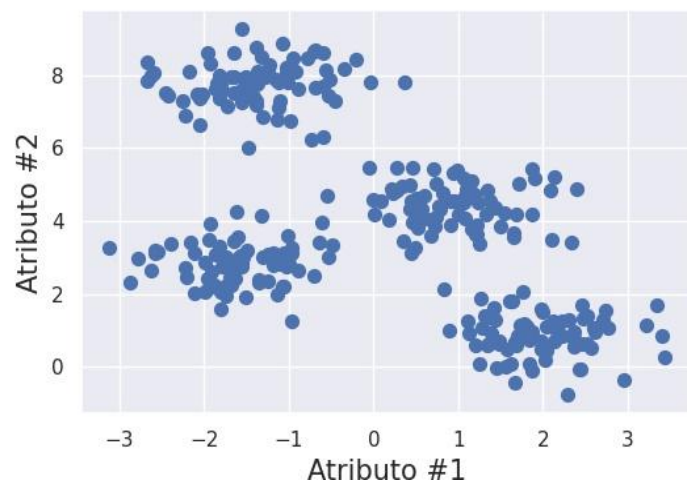


Agrupamento Hierárquico

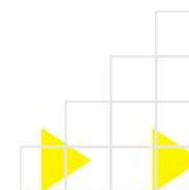
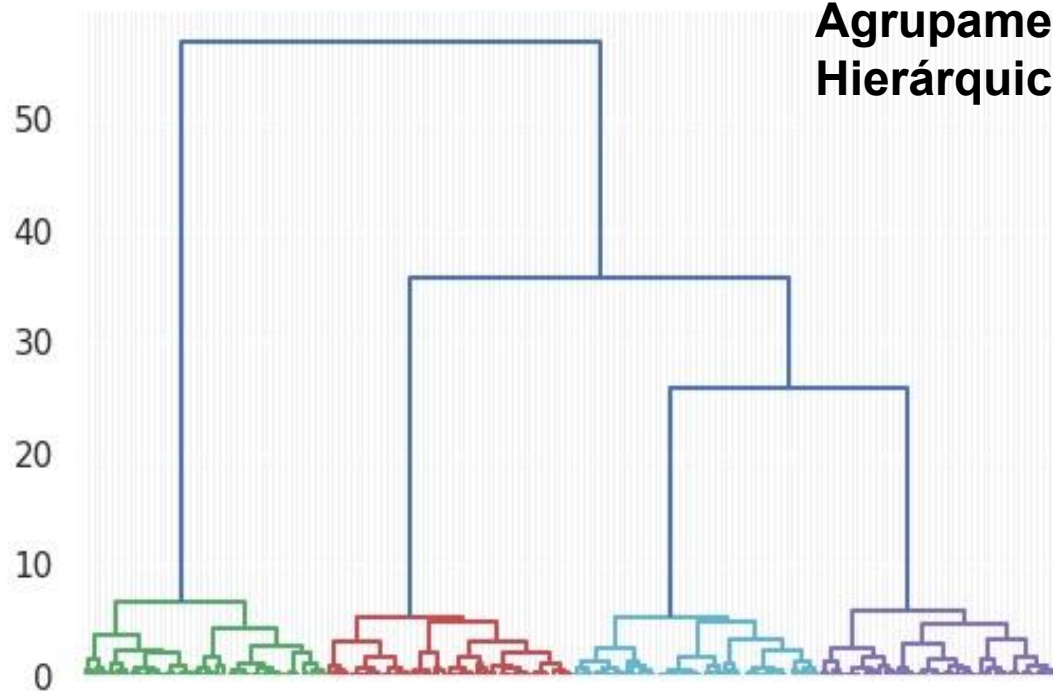


- Podemos inspecionar o dendrograma para estimar o número natural de clusters. No exemplo, há 4 subárvores bem separadas.

Conjunto de Dados



Agrupamento Hierárquico

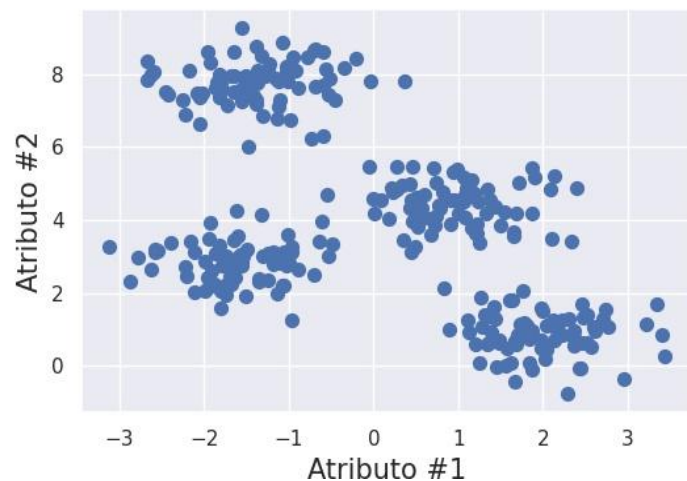


Agrupamento Hierárquico

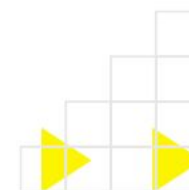
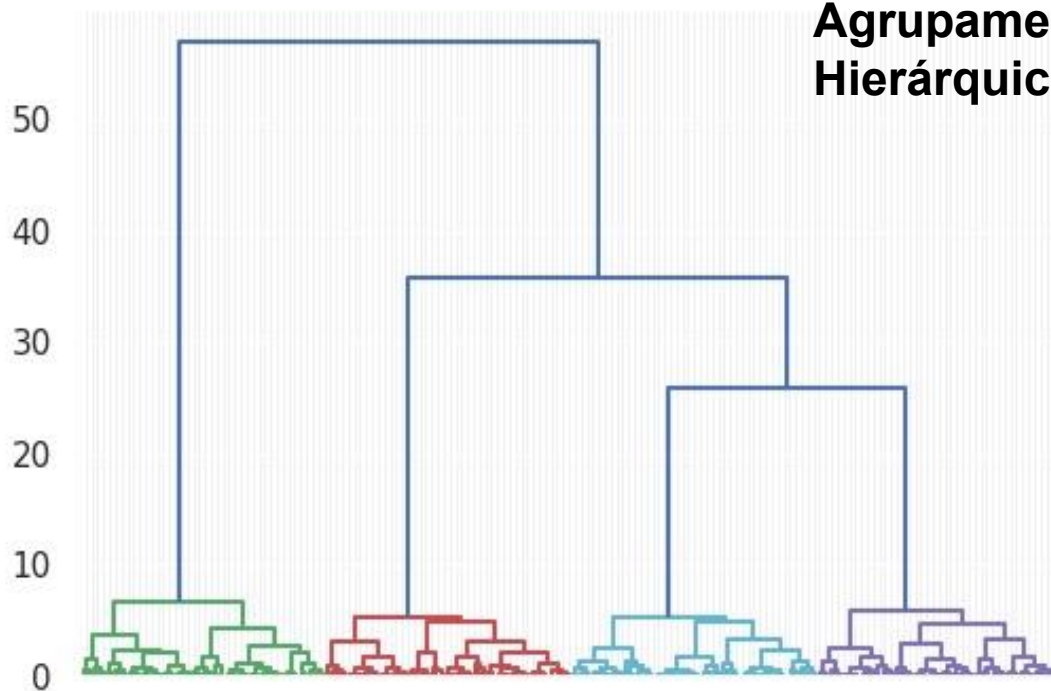


- Conceitos de homogeneidade (coesão interna) e heterogeneidade (separabilidade) dos clusters representados pela altura (eixo y) da união entre cluster.

Conjunto de Dados



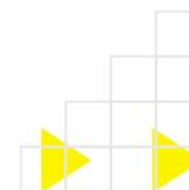
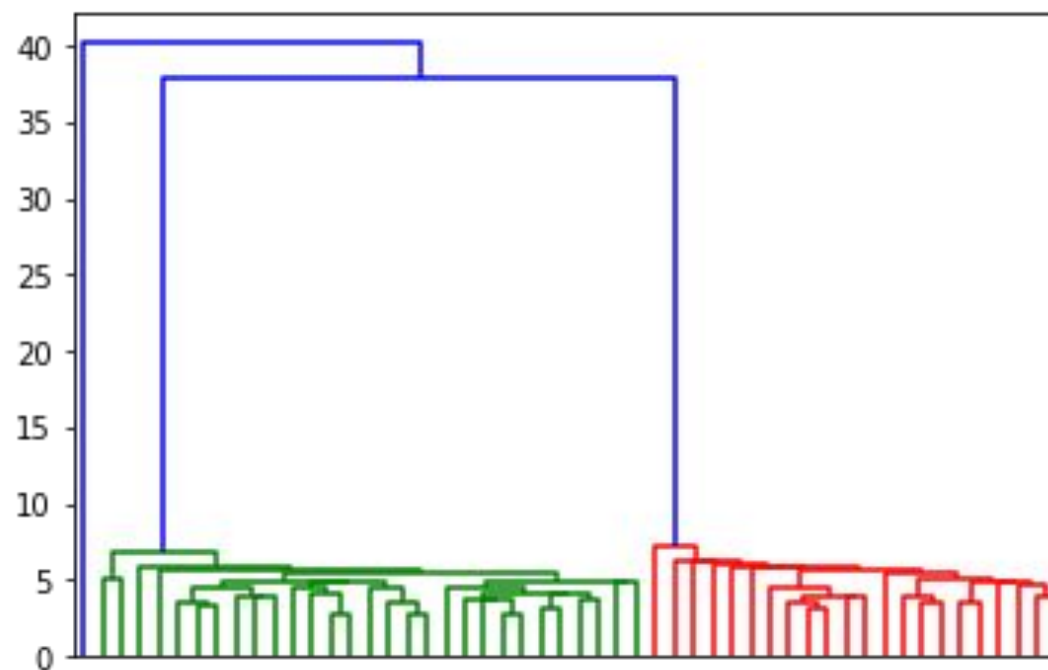
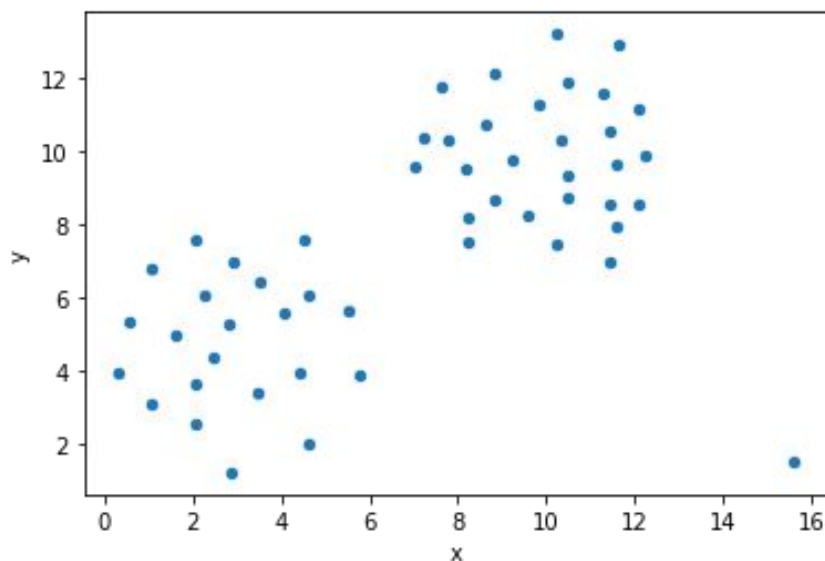
Agrupamento Hierárquico



Agrupamento Hierárquico



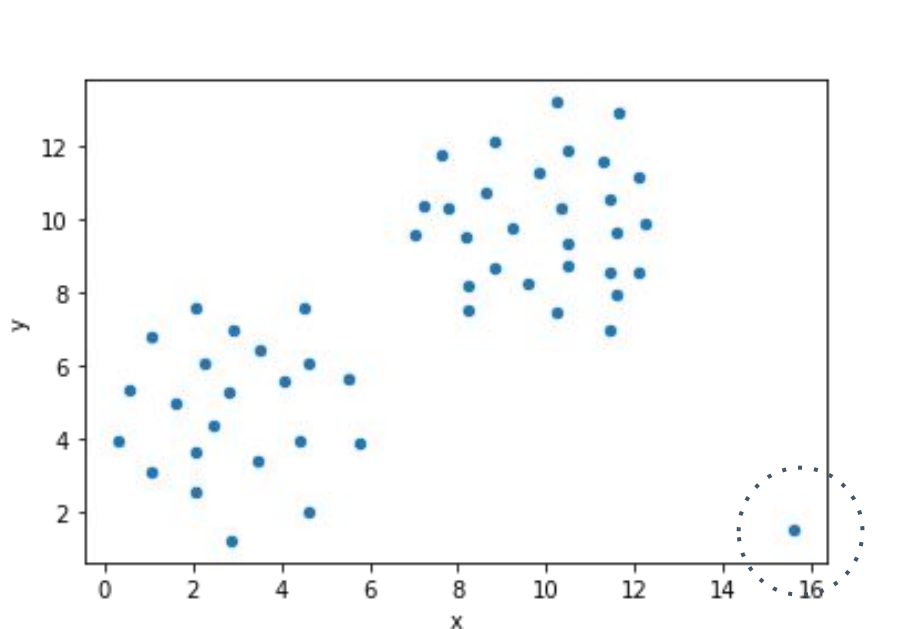
- Dendrogramas podem ser úteis para detectar *outliers*.
- Ramos isolados indicam objetos muito distante dos demais.



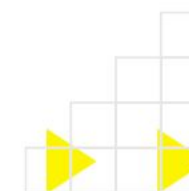
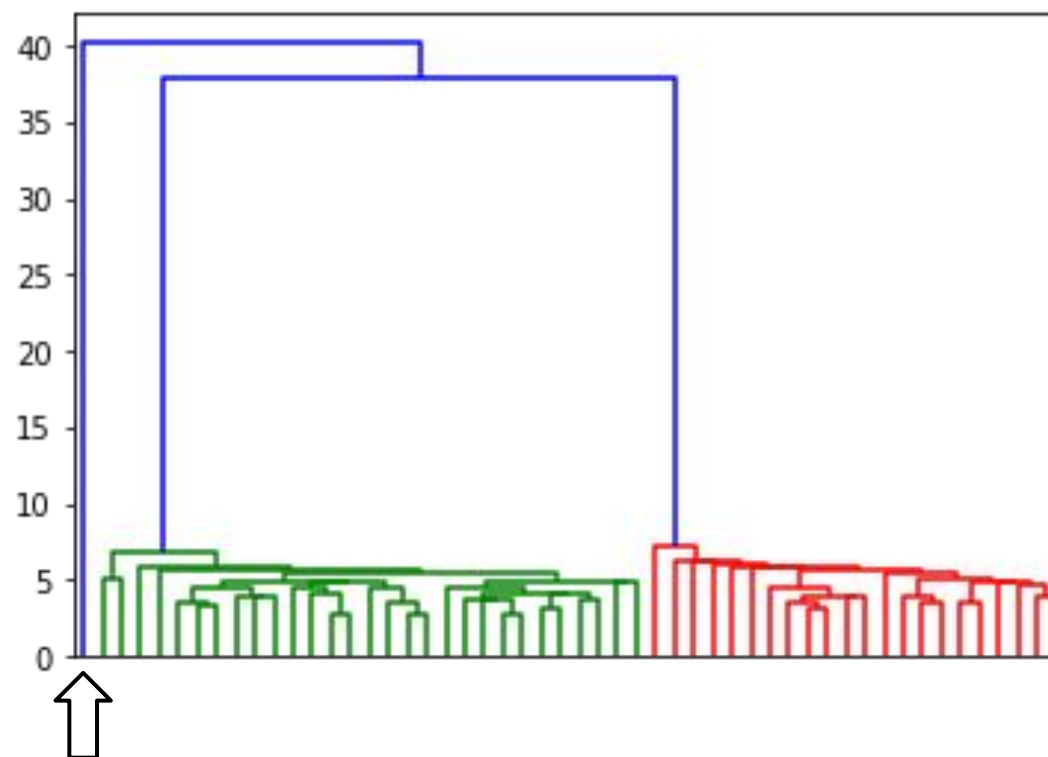
Agrupamento Hierárquico



- Dendrogramas podem ser úteis para detectar *outliers*.
- Ramos isolados indicam objetos muito distante dos demais.



Outlier



Agrupamento Hierárquico



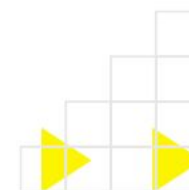
- Dois métodos clássicos para agrupamento hierárquico

Aglomerativos:

- Iniciar alocando cada objeto em um *cluster*
- Encontrar o melhor par de *clusters* para unir
- Repetir até formar um único *cluster*

Divisivos:

- Iniciar alocando todos os objetos em um único *cluster*
- Dividir um *cluster* em dois *subclusters*
- Repetir a divisão até que cada objeto seja um *cluster*



Agrupamento Hierárquico



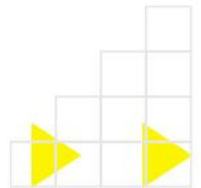
- Dois métodos clássicos para agrupamento hierárquico

Aglomerativos:

- Iniciar alocando cada objeto em um *cluster*
- Encontrar o melhor par de *clusters* para unir
- Repetir até formar um único *cluster*

Divisivos:

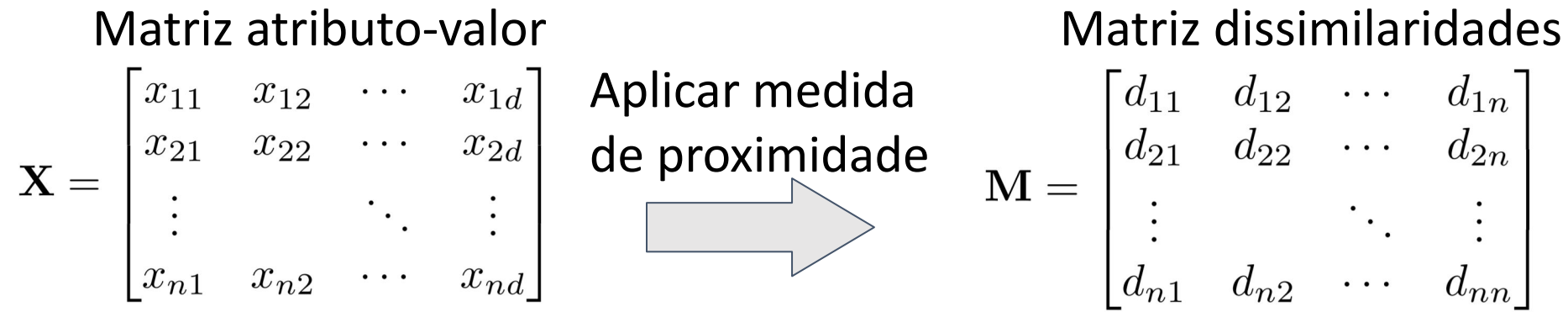
- Iniciar alocando todos os objetos em um único *cluster*
- Dividir um *cluster* em dois *subclusters*
- Repetir a divisão até que cada objeto seja um *cluster*



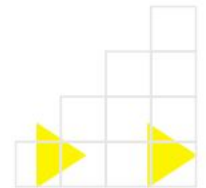
Agrupamento Hierárquico



- **Matriz de dissimilaridades:** armazena as distâncias entre cada par de objetos do conjunto de dados



- As diagonais são dissimilaridades entre um mesmo objeto
- Por simetria, podemos usar apenas a triangular superior ou inferior da matriz de dissimilaridades



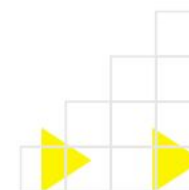
Agrupamento Hierárquico



- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos

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



Agrupamento Hierárquico

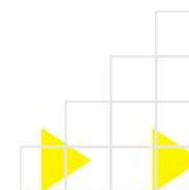


- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos

Inicialmente, cada objeto é um cluster.

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



Agrupamento Hierárquico

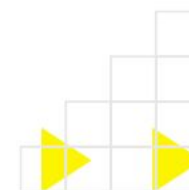


- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

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



Agrupamento Hierárquico

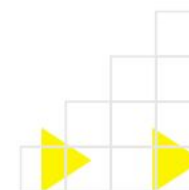
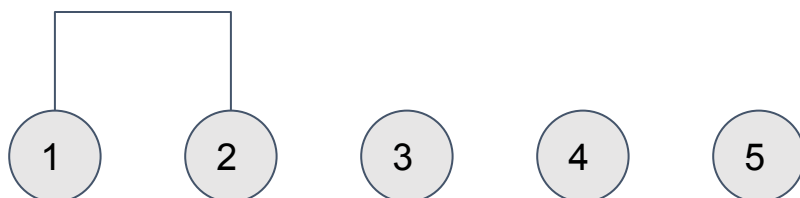


- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

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



Agrupamento Hierárquico

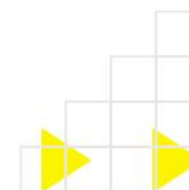
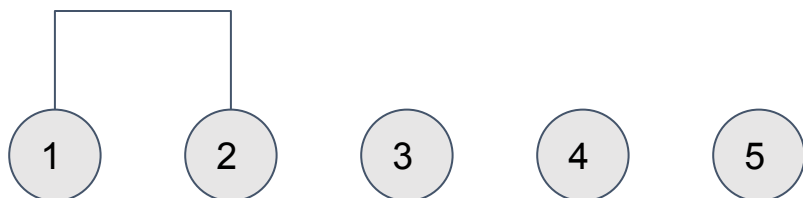


- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos

Atualizar a matriz de dissimilaridades!

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



Agrupamento Hierárquico



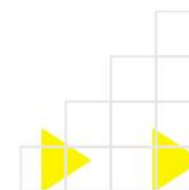
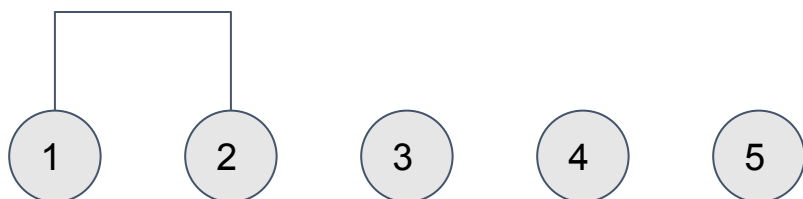
- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

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

Atualizar a matriz de dissimilaridades!

	12	3	4	5
12	0			
3		0		
4		4	0	
5		5	3	0



Agrupamento Hierárquico



- Exemplo de Agrupamento Hierárquico

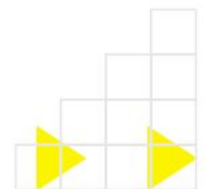
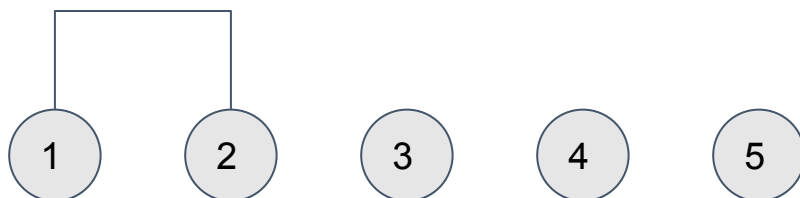
Considere uma matriz de dissimilaridades calculada para 5 objetos

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

Atualizar a matriz de dissimilaridades!

$$\begin{aligned}d_{(12)3} &= \min \{d_{13}, d_{23}\} = d_{23} = 5 \\d_{(12)4} &= \min \{d_{14}, d_{24}\} = d_{24} = 9 \\d_{(12)5} &= \min \{d_{15}, d_{25}\} = d_{25} = 8\end{aligned}$$

	12	3	4	5
12	0			
3		0		
4		4	0	
5		5	3	0



Agrupamento Hierárquico



- Exemplo de Agrupamento Hierárquico

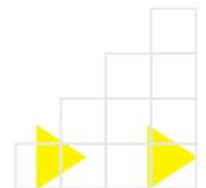
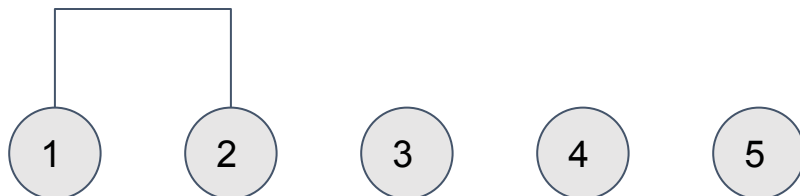
Considere uma matriz de dissimilaridades calculada para 5 objetos

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

Atualizar a matriz de dissimilaridades!

$$\begin{aligned}d_{(12)3} &= \min \{d_{13}, d_{23}\} = d_{23} = 5 \\d_{(12)4} &= \min \{d_{14}, d_{24}\} = d_{24} = 9 \\d_{(12)5} &= \min \{d_{15}, d_{25}\} = d_{25} = 8\end{aligned}$$

	12	3	4	5
12	0			
3	5	0		
4	9	4	0	
5	8	5	3	0



Agrupamento Hierárquico

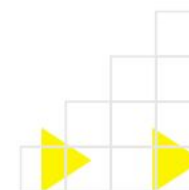
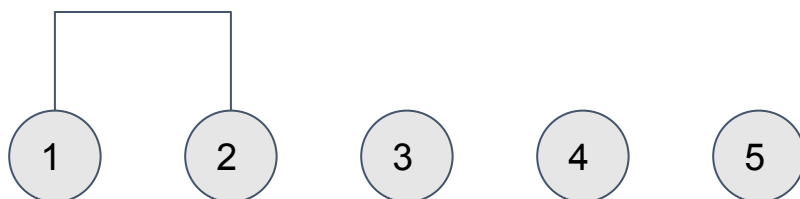


- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

	12	3	4	5
12	0			
3	5	0		
4	9	4	0	
5	8	5	3	0



Agrupamento Hierárquico

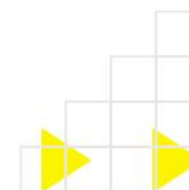
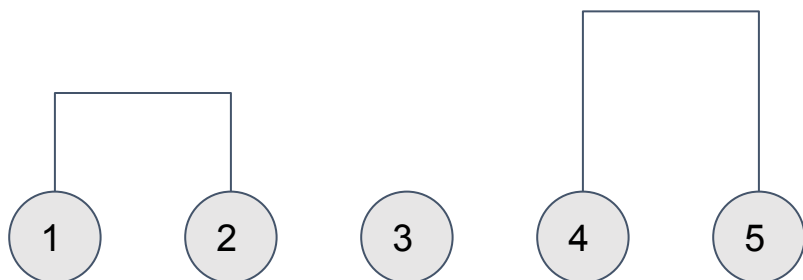


- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

	12	3	4	5
12	0			
3	5	0		
4	9	4	0	
5	8	5	3	0



Agrupamento Hierárquico



- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

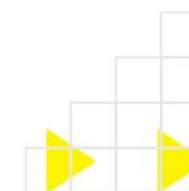
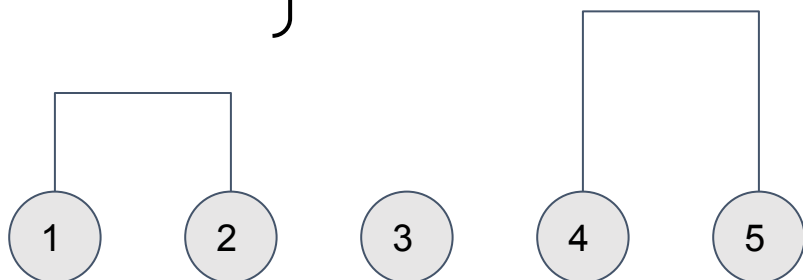
	1	2	3	4	5
1	0				
2	2	0			
3	6	5	0		
4	10	9	4	0	
5	9	8	5	3	0

Atualizar a matriz de dissimilaridades!

$$d_{(12)(45)} = \min \{d_{14}, d_{15}, d_{24}, d_{25}\} = d_{25} = 8$$

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

	12	3	45
12	0		
3	5	0	
45			0



Agrupamento Hierárquico



- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

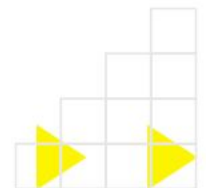
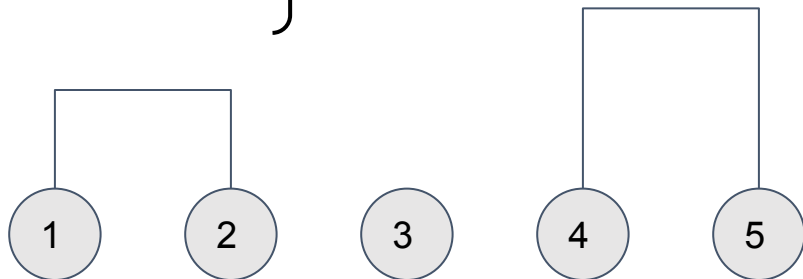
	1	2	3	4	5
1	0				
2	2	0			
3	6	5	0		
4	10	9	4	0	
5	9	8	5	3	0

Atualizar a matriz de dissimilaridades!

$$d_{(12)(45)} = \min \{d_{14}, d_{15}, d_{24}, d_{25}\} = d_{25} = 8$$

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

	12	3	45
12	0		
3	5	0	
45	8	4	0



Agrupamento Hierárquico

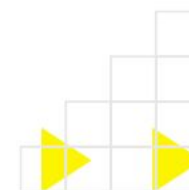
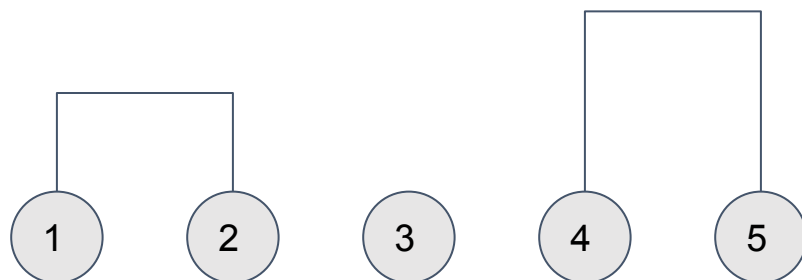


- Exemplo de Agrupamento Hierárquico

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

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



Agrupamento Hierárquico

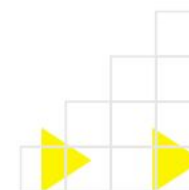
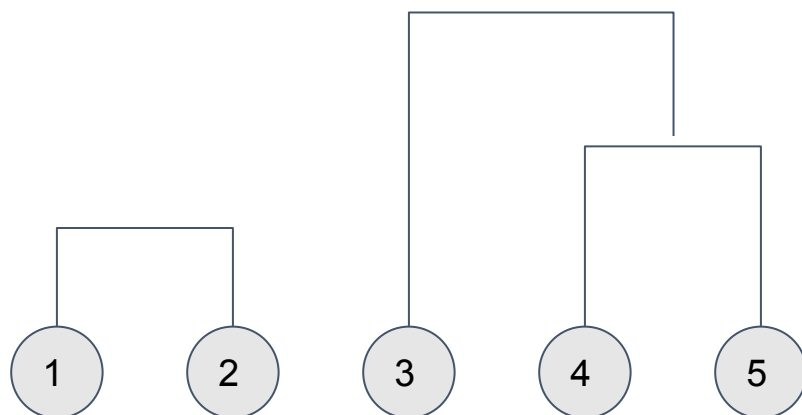


- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos

Qual é o melhor par de *clusters* para unir?

$$\begin{matrix} & \mathbf{12} & \mathbf{3} & \mathbf{45} \\ \mathbf{12} & \begin{bmatrix} 0 & & \\ 5 & 0 & \\ 8 & \textcircled{4} & 0 \end{bmatrix} \\ \mathbf{3} & & & \\ \mathbf{45} & & & \end{matrix}$$



Agrupamento Hierárquico

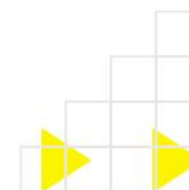
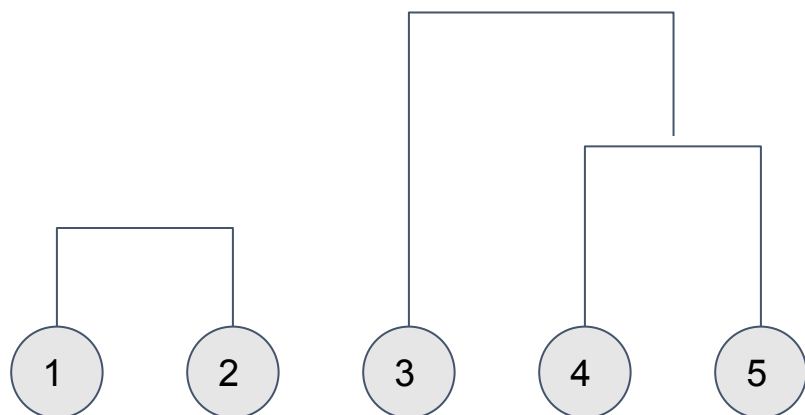


- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos

Atualizar a matriz de dissimilaridades!

$$\begin{matrix} & \mathbf{12} & \mathbf{345} \\ \mathbf{12} & \begin{bmatrix} 0 & & \end{bmatrix} \\ \mathbf{345} & \begin{bmatrix} 5 & 0 \end{bmatrix} \end{matrix}$$

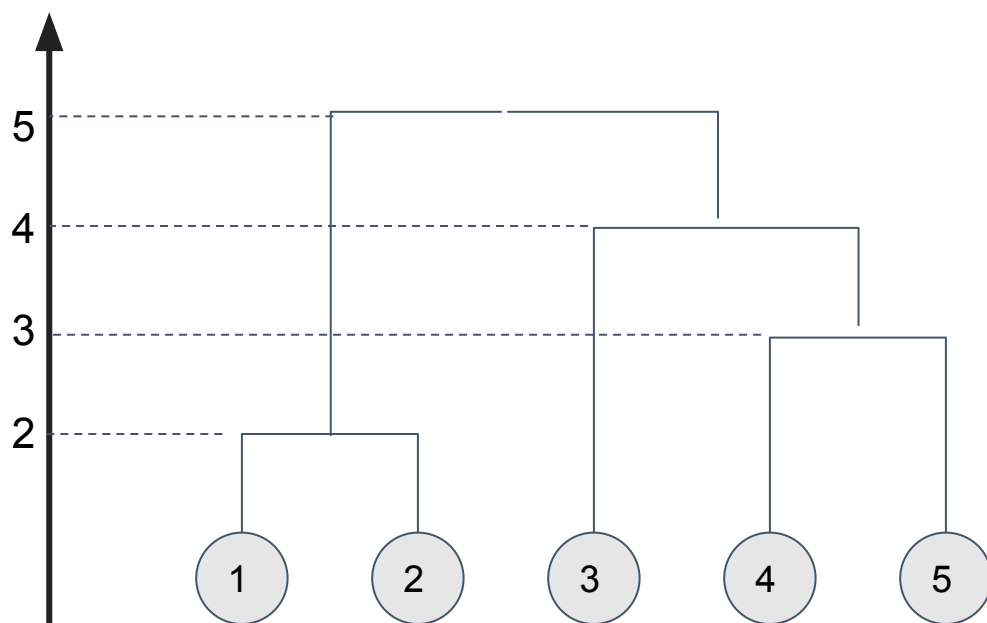


Agrupamento Hierárquico

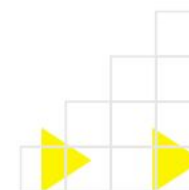


- **Exemplo de Agrupamento Hierárquico**

Considere uma matriz de dissimilaridades calculada para 5 objetos



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

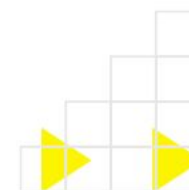
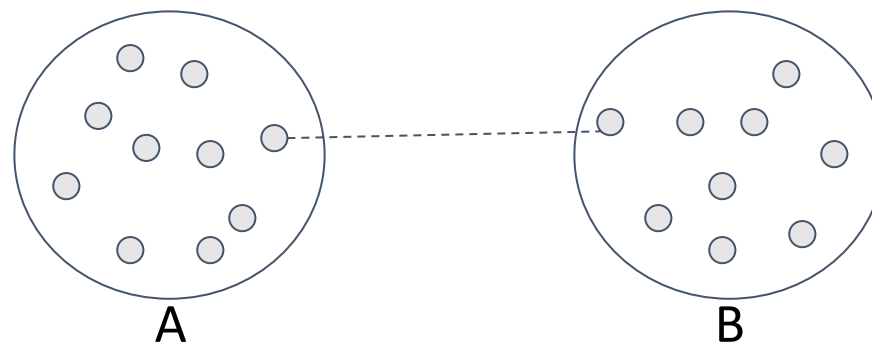


Agrupamento Hierárquico



- A etapa mais importante é determinar distância entre clusters!

Nosso exemplo utilizou a estratégia de Vizinho mais Próximo (MIN)



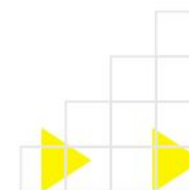
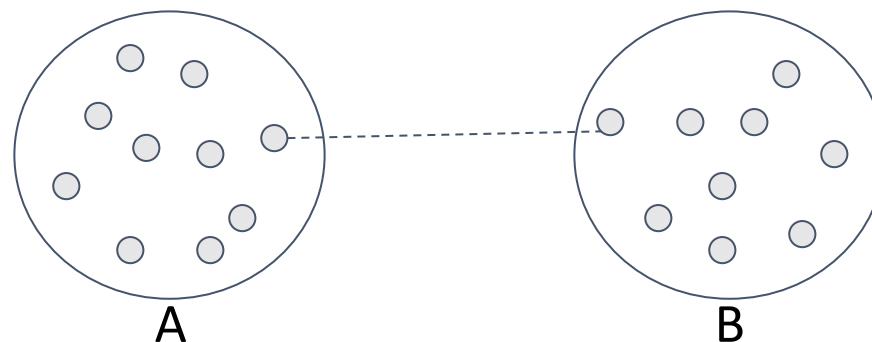
Agrupamento Hierárquico



- A etapa mais importante é determinar distância entre clusters!

Nosso exemplo utilizou a estratégia de Vizinho mais Próximo (MIN)

Distância entre clusters A e B = menor distância de qualquer objeto do cluster A para qualquer objeto do cluster B



Agrupamento Hierárquico

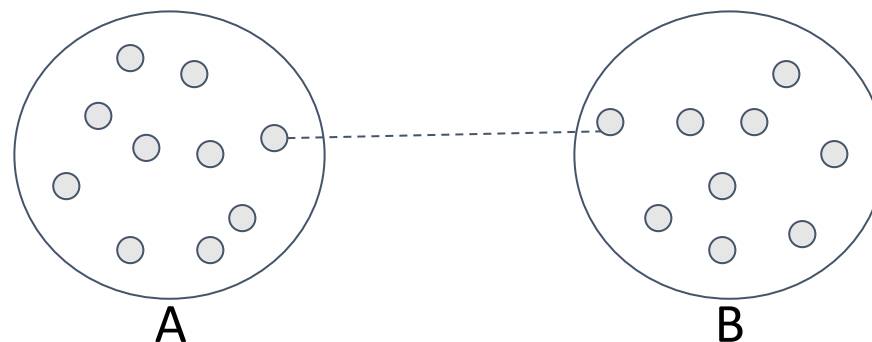


- A etapa mais importante é determinar distância entre clusters!

Nosso exemplo utilizou a estratégia de Vizinho mais Próximo (MIN)

Distância entre clusters A e B = menor distância de qualquer objeto do cluster A para qualquer objeto do cluster B

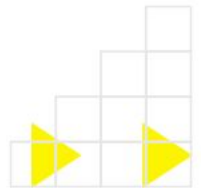
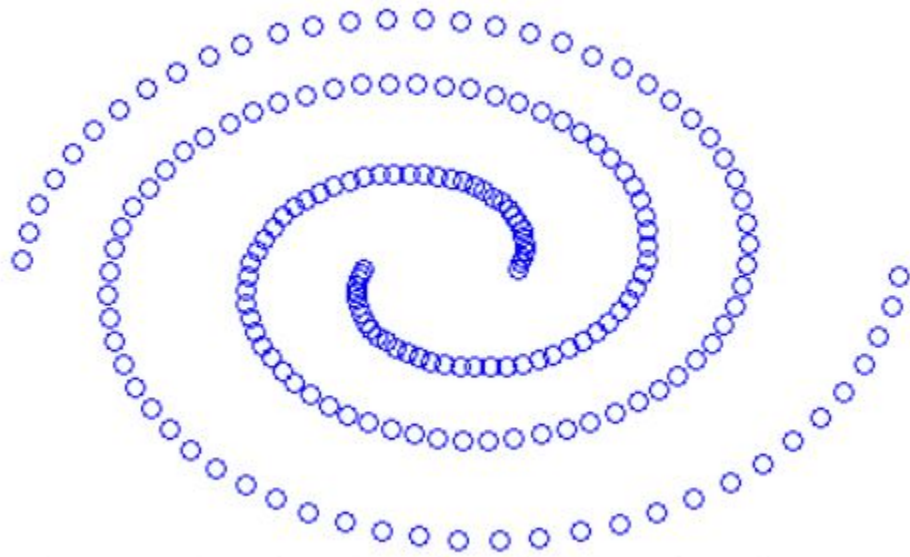
Single-Link



Agrupamento Hierárquico

- Single-Link (MIN)

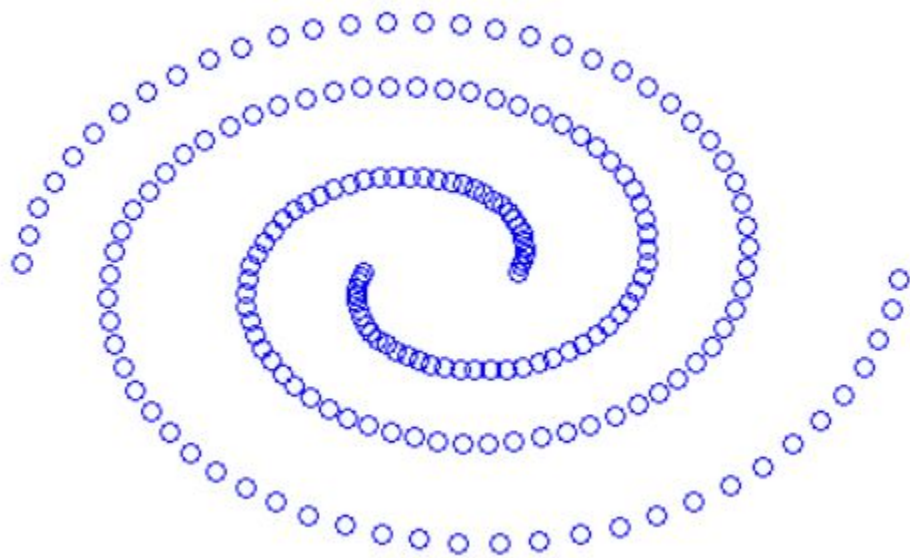
Tende a capturar *clusters* baseado em contiguidade ou encadeamento.



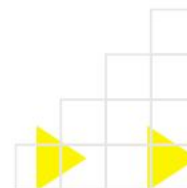
Agrupamento Hierárquico

- Single-Link (MIN)

Tende a capturar *clusters* baseado em contiguidade ou encadeamento.

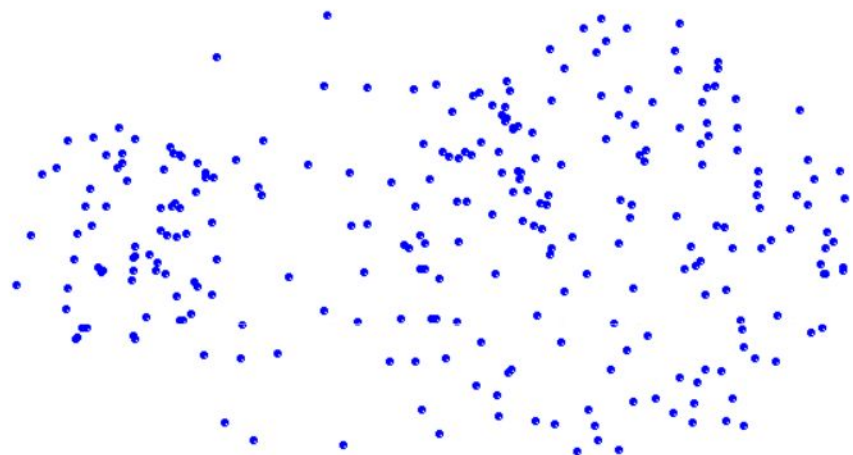


Sensível a ruídos!

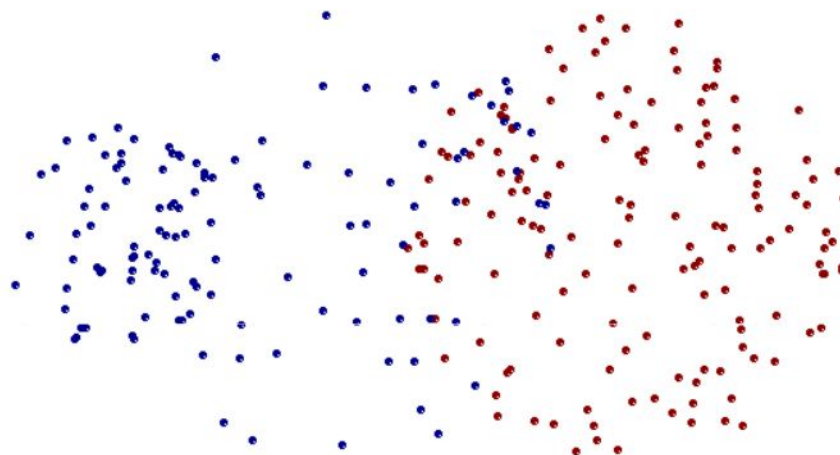


Agrupamento Hierárquico

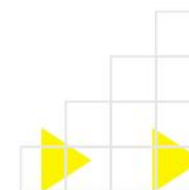
- Single-Link (MIN)



Pontos Originais



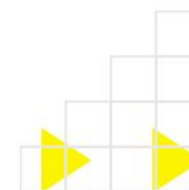
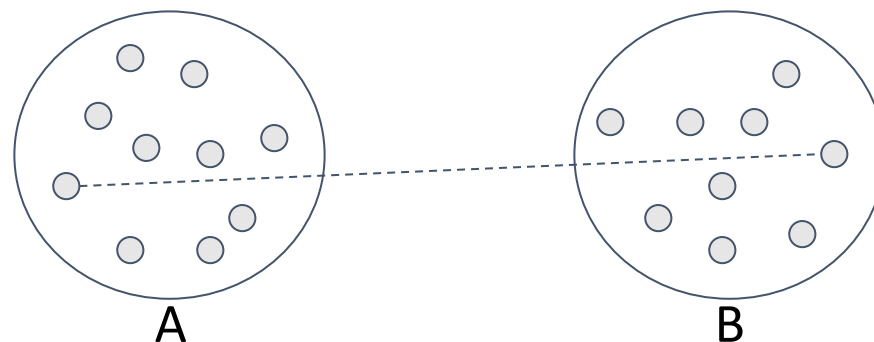
Pontos Agrupados
(dois clusters)



Agrupamento Hierárquico

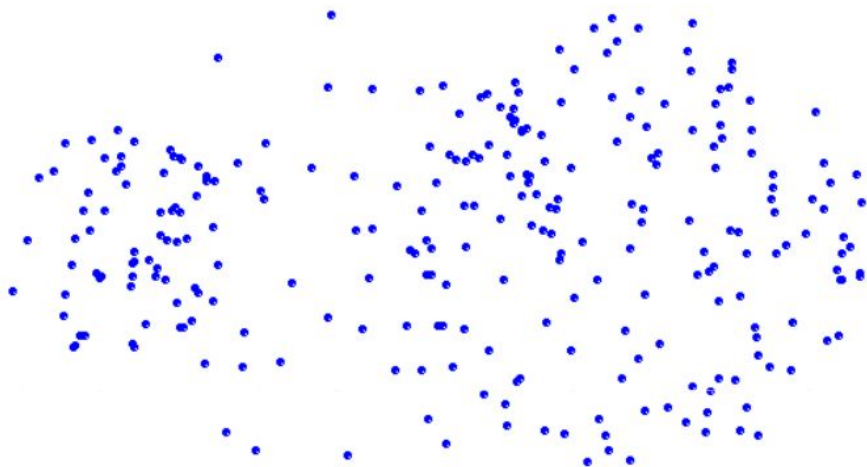
- Complete-Link (MAX)

Distância entre clusters A e B = maior distância de qualquer objeto do cluster A para qualquer objeto do cluster B

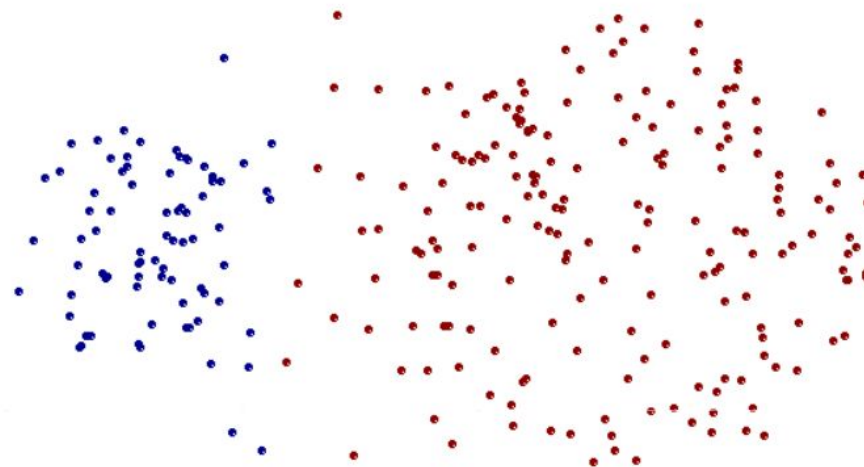


Agrupamento Hierárquico

- Complete-Link (MAX)

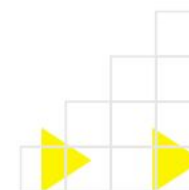


Pontos Originais



Pontos Agrupados
(dois clusters)

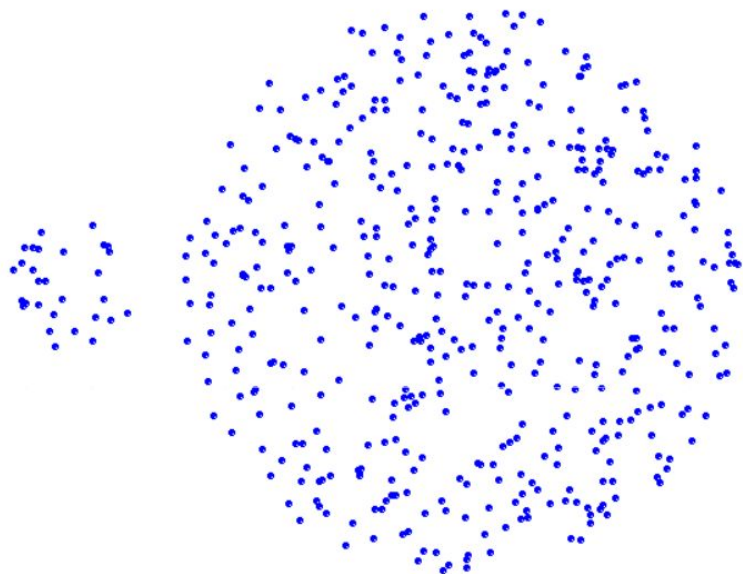
Tende a encontrar clusters com diâmetros semelhantes.
Vantagem: menos sensível a ruído.



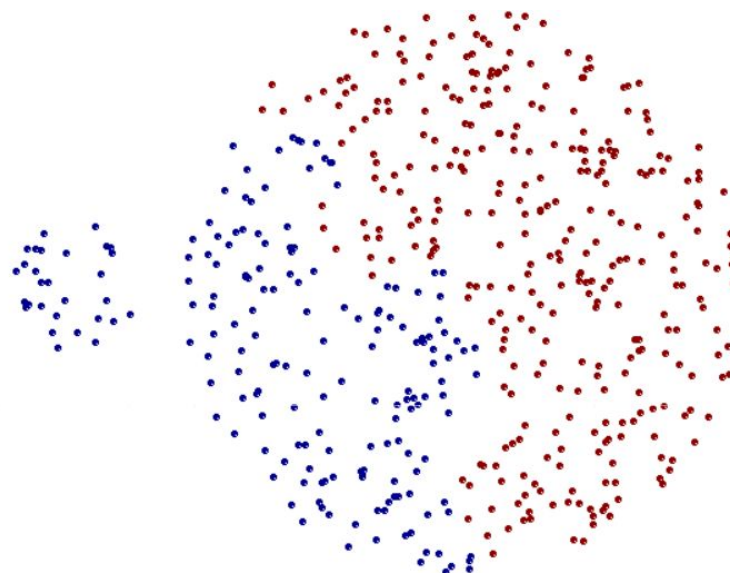
Agrupamento Hierárquico



- Complete-Link (MAX)

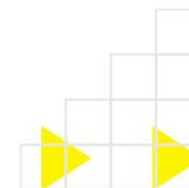


Pontos Originais



Pontos Agrupados
(dois clusters)

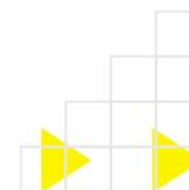
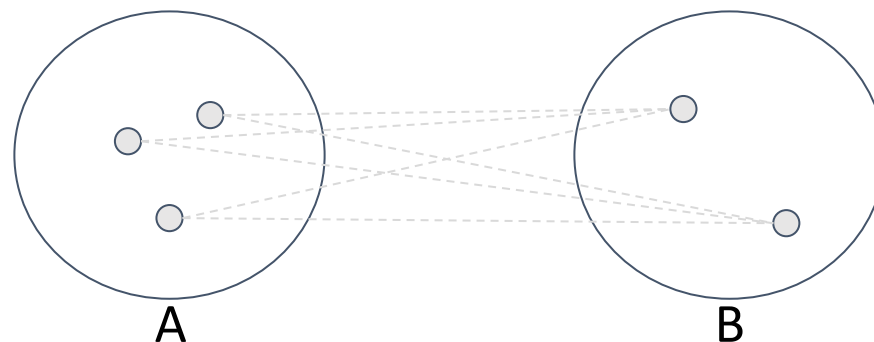
Tende a encontrar clusters com diâmetros semelhantes.
Desvantagem: tende a dividir grandes *clusters*.



Agrupamento Hierárquico

- Average-Link (Média)

Distância entre clusters A e B = média das distâncias de todos os objetos do *cluster A* em relação aos objetos do *cluster B*

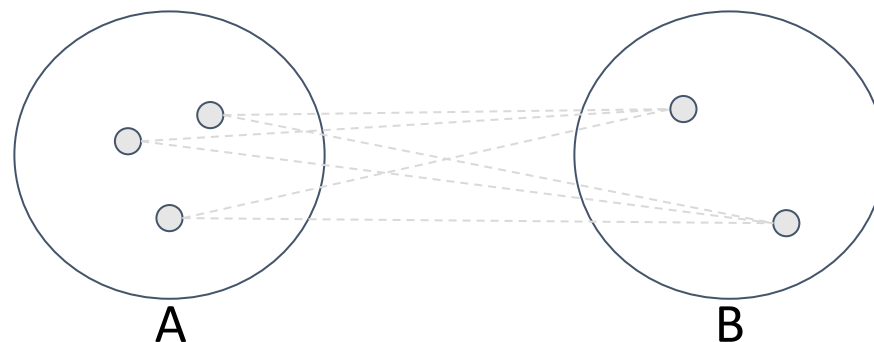


Agrupamento Hierárquico

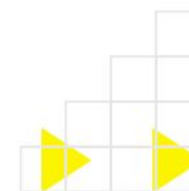


- **Average-Link (Média)**

Distância entre clusters A e B = média das distâncias de todos os objetos do *cluster* A em relação aos objetos do *cluster* B



$$d(C_i, C_j) = \frac{\sum_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)}{|C_i||C_j|}$$



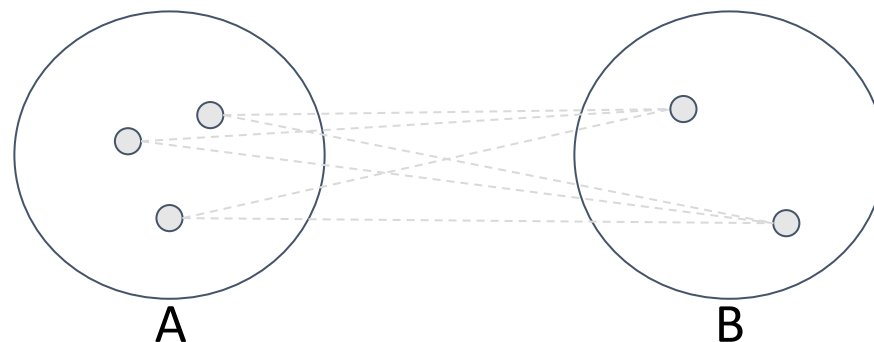
Agrupamento Hierárquico



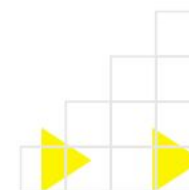
- **Average-Link (Média)**

Distância entre clusters A e B = média das distâncias de todos os objetos do *cluster A* em relação aos objetos do *cluster B*

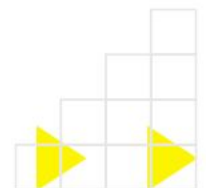
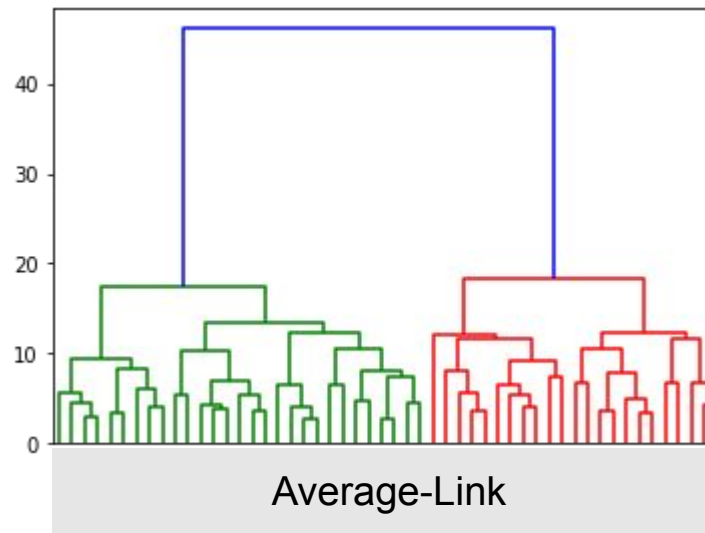
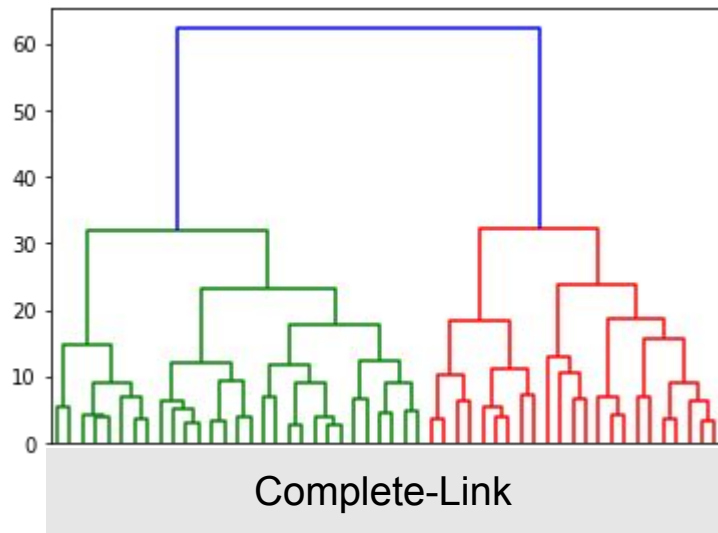
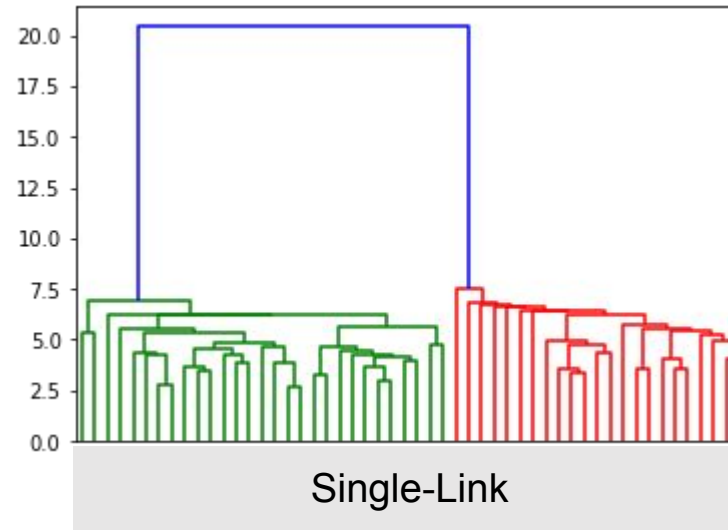
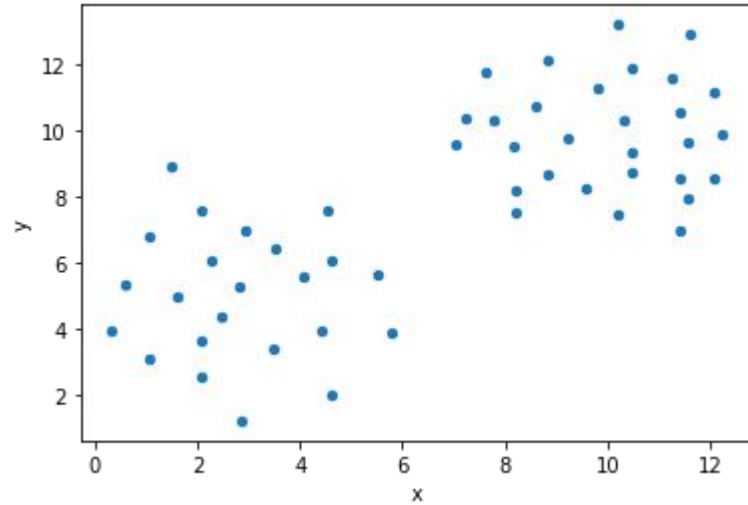
Ponto de
equilíbrio entre
Single-Link e
Complete-Link



$$d(C_i, C_j) = \frac{\sum_{\mathbf{x}_i \in C_i, \mathbf{x}_j \in C_j} d(\mathbf{x}_i, \mathbf{x}_j)}{|C_i||C_j|}$$



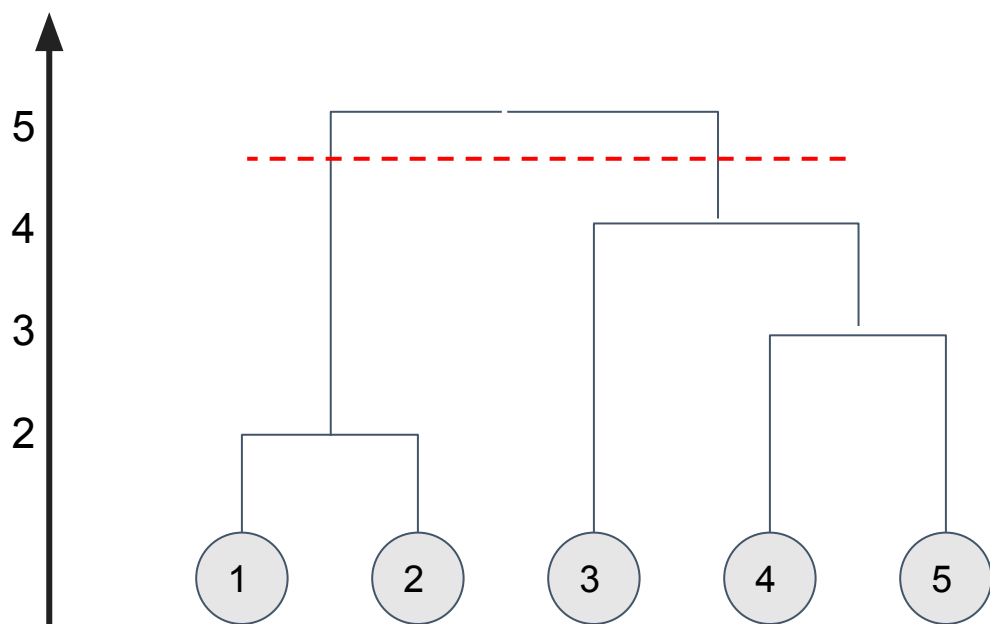
Agrupamento Hierárquico



Agrupamento Hierárquico

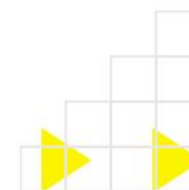


- Podemos extrair partições a partir de um agrupamento hierárquico
- Escolha do número de clusters “*a posteriori*”



Cortes no dendrograma

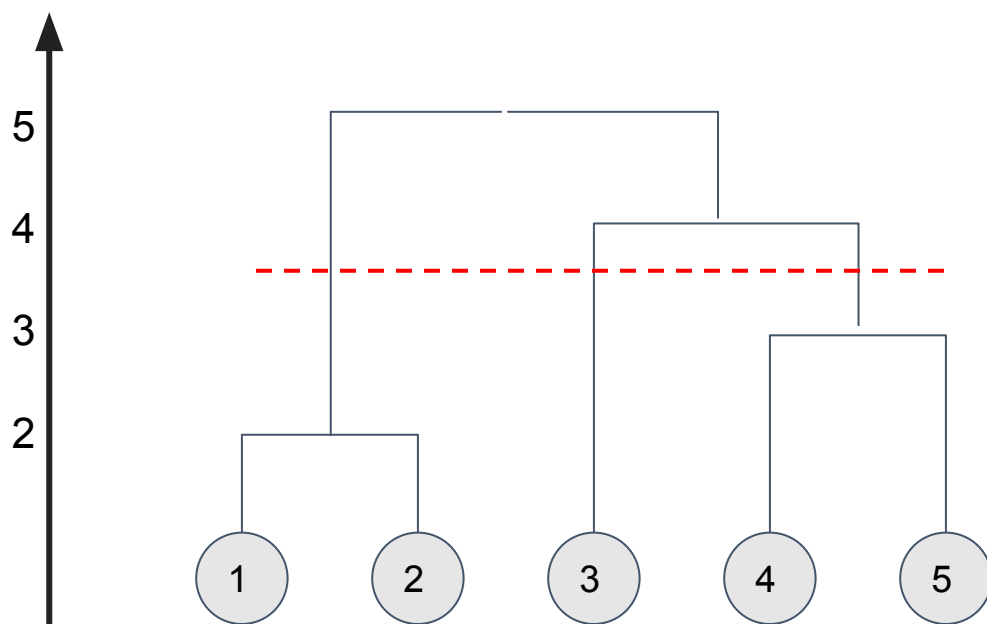
Partição $P = \{(1,2), (3,4,5)\}$



Agrupamento Hierárquico

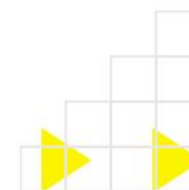


- Podemos extrair partições a partir de um agrupamento hierárquico
- Escolha do número de clusters “*a posteriori*”



Cortes no dendrograma

Partição $P = \{(1,2), (3), (4,5)\}$



Bibliografia

Rezende, S. O. (2003). *Sistemas inteligentes: fundamentos e aplicações*. Editora Manole Ltda.

Tan, P.N.; Steinbach, M.; Karpatne, A.; Kumar, V. (2016). *Introduction to Data Mining (2nd Edition)*. Pearson.

