

MBA⁺

**Artificial Intelligence &
Machine Learning**





Aprendizado não supervisionado

K-means / EM / Regras de Associação



Algoritmos Hierárquicos

Hierarquia: Conceitos Básicos



Hierarquias são comumente usadas para organizar informação

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](#)...

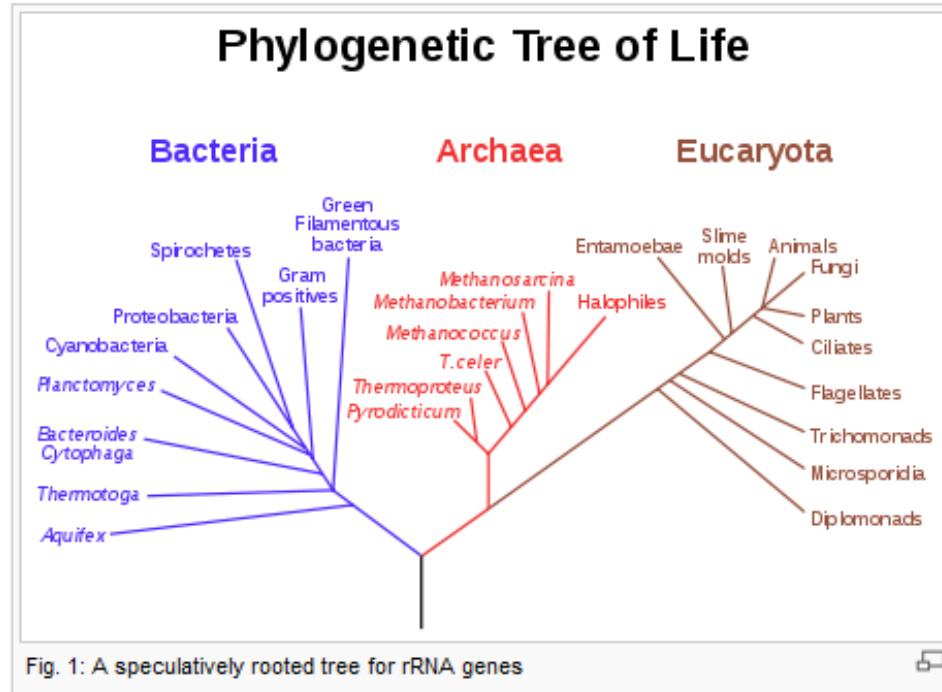
Society & Culture

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



Hierarquia: Conceitos Básicos

Exemplo: árvores filogenéticas em biologia

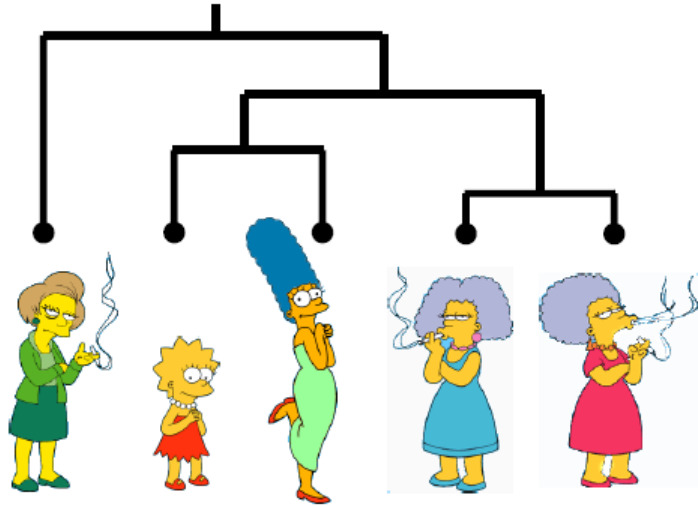


Métodos Clássicos para Agrupamento Hierárquico



Bottom-Up (aglomerativos):

- Iniciar colocando cada objeto em um *cluster*
- Encontrar o melhor par de *clusters* para unir
- Unir o par de *clusters* escolhido
- Repetir até que todos os objetos estejam reunidos em um só *cluster*



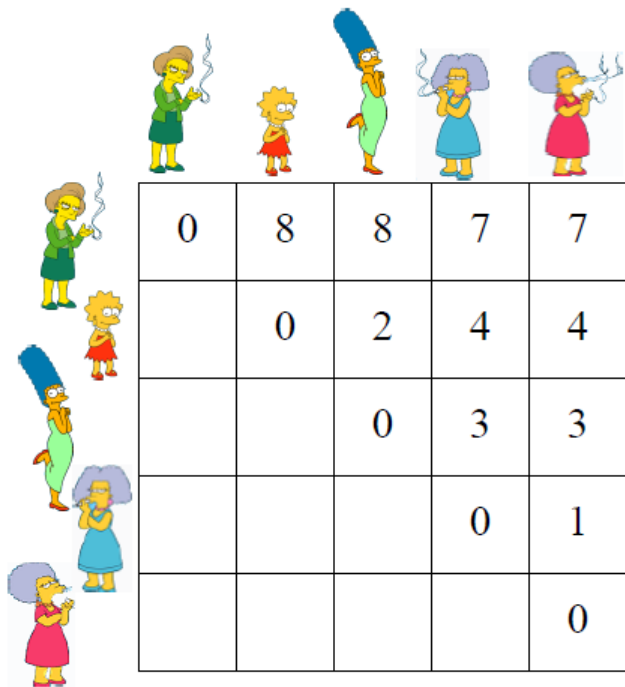
Top-Down (divisivos):

- 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

Métodos Clássicos para Agrupamento Hierárquico



Algoritmos hierárquicos podem operar somente sobre uma matriz de distâncias.



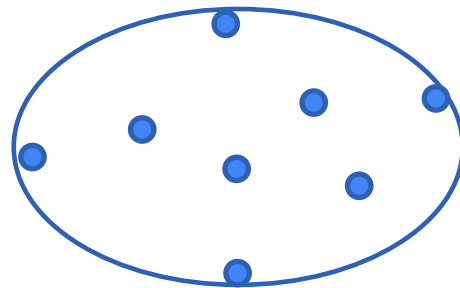
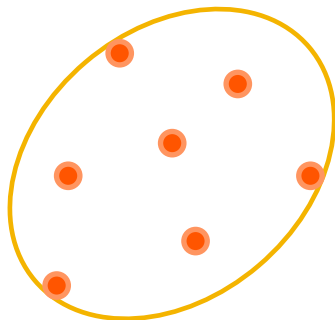
$$D(\text{Maggie}, \text{Patty}) = 1$$
$$D(\text{Mrs. Simpson}, \text{Lisa}) = 8$$

Como definir Inter-Cluster (Dis)similaridade



Matriz de Distância

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



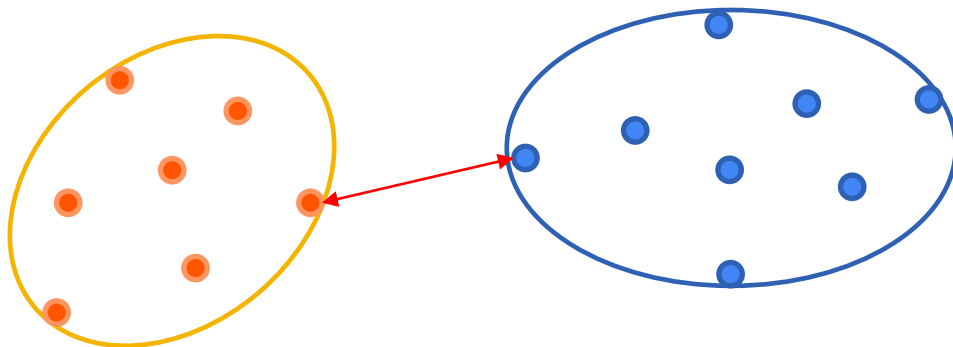
- MIN
- MAX
- Group Average
- ...

Como definir Inter-Cluster (Dis)similaridade



Matriz de Distância

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

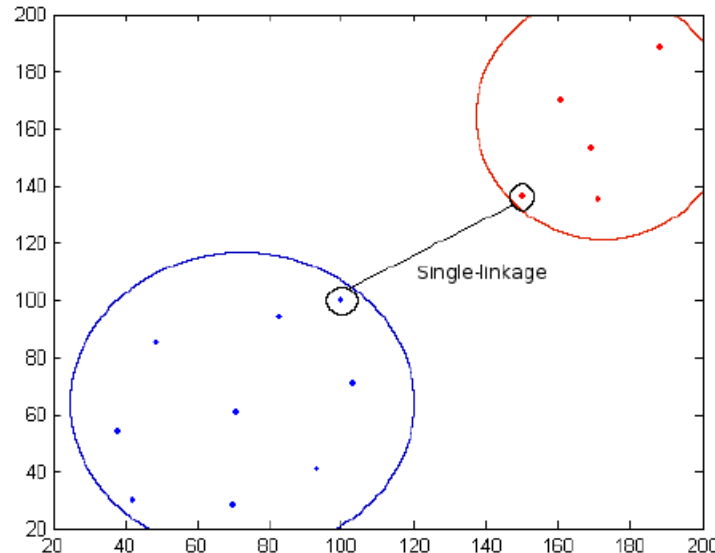


- MIN
- MAX
- Group Average
- ...

Single Linkage (Florek, 1951)



- Dissimilaridade entre clusters é dada pela **menor** dissimilaridade entre 2 objetos (um de cada cluster)
 - Originalmente baseado em Grafos: **menor** aresta entre dois vértices de subconjuntos distintos

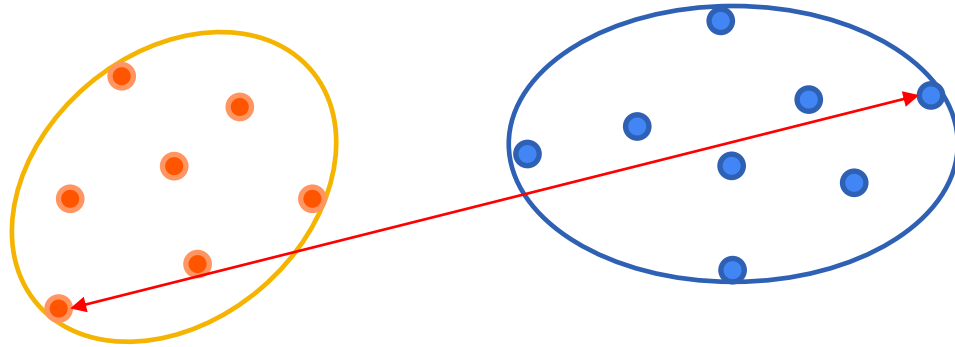


Como definir Inter-Cluster (Dis)similaridade



Matriz de Distância

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

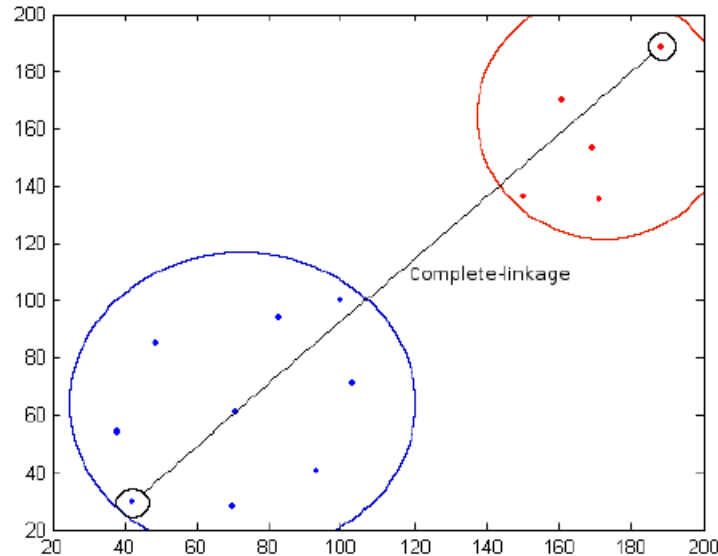


- MIN
- **MAX**
- Group Average
- ...

Complete Linkage (Sorensen, 1948)



- Dissimilaridade entre clusters é dada pela **maior** dissimilaridade entre 2 objetos (um de cada cluster)
 - Originalmente baseado em Grafos: maior aresta entre dois vértices de subconjuntos distintos

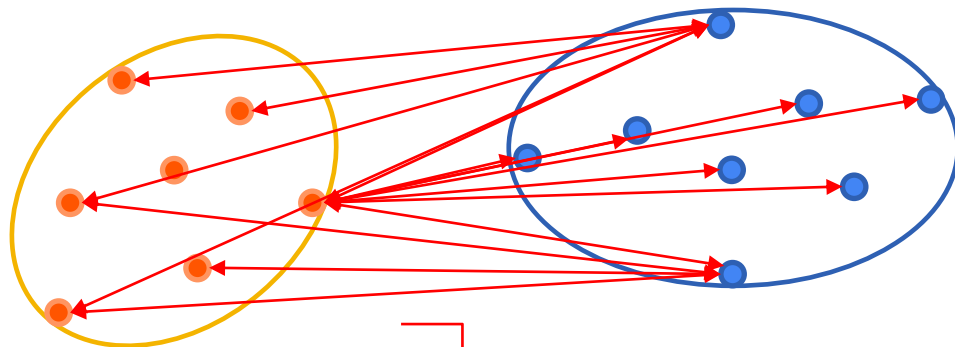


Como definir Inter-Cluster (Dis)similaridade



Matriz de Distância

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



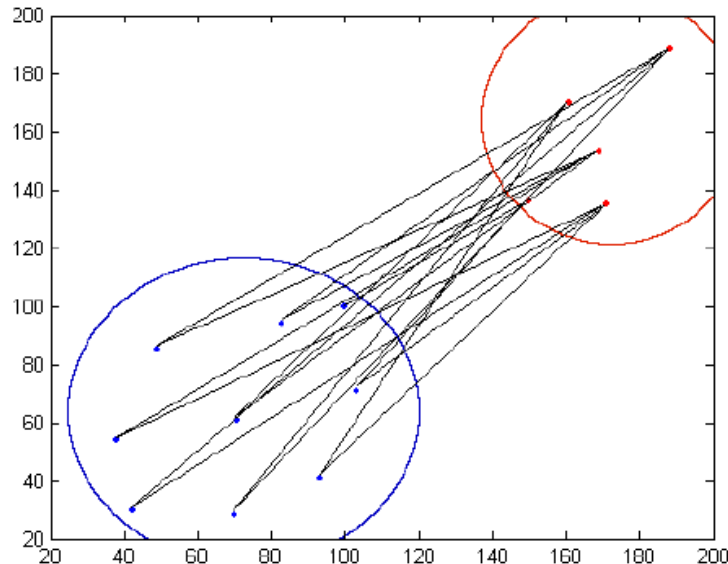
Todas as distâncias para-a-par

- MIN
- MAX
- Group Average
- ...

Complete Linkage (Sokal R and Michener C, 1958)



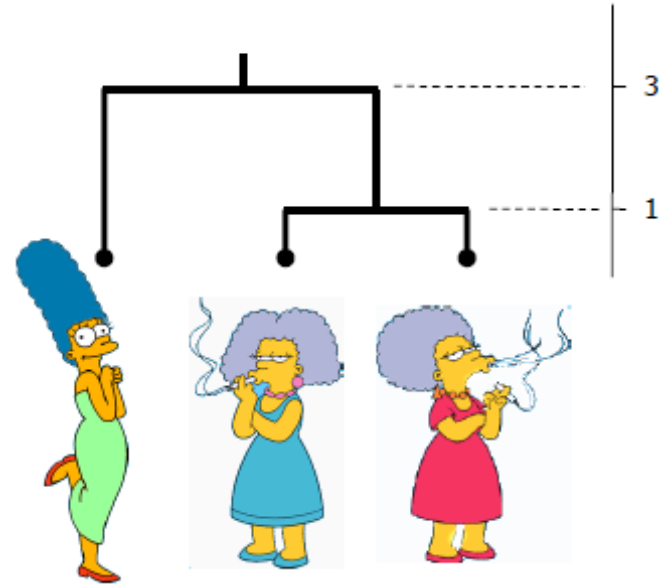
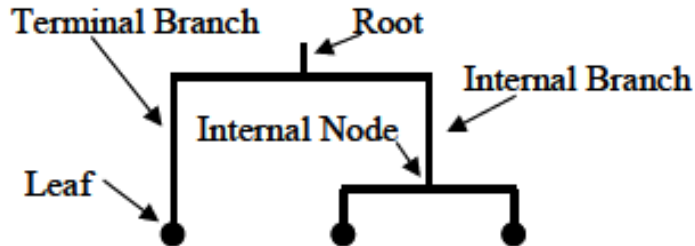
- Dissimilaridade entre clusters é dada pela **distância média** entre cada par de objetos (um de cada cluster)
- Também conhecido como UPGMA – Unweighted Pair Group Method using Arithmetic averages



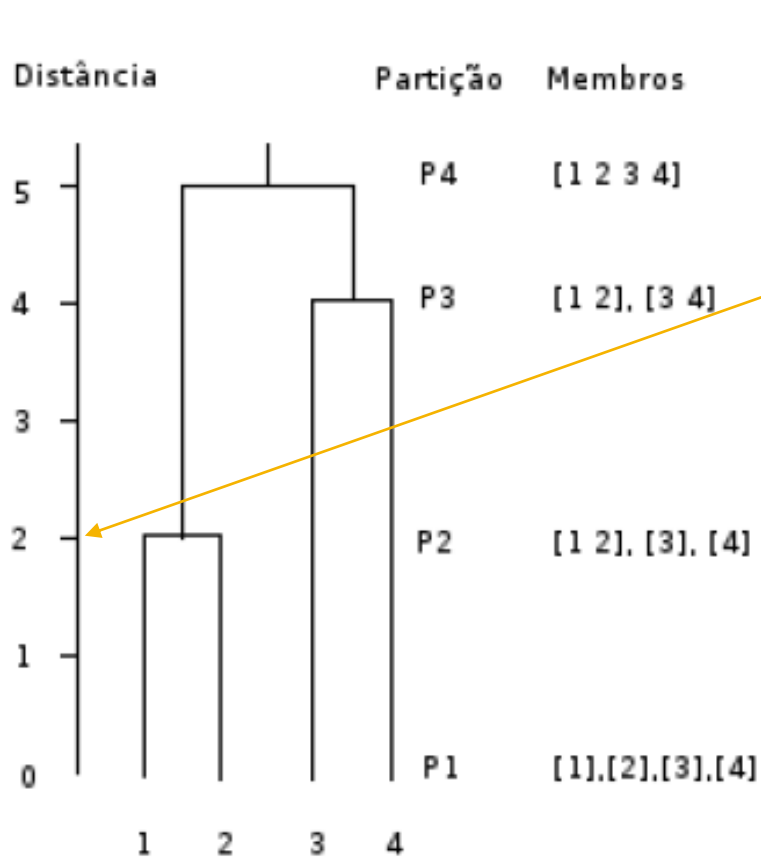
Dendrograma = Hierarquia + Dissimilaridade entre Clusters



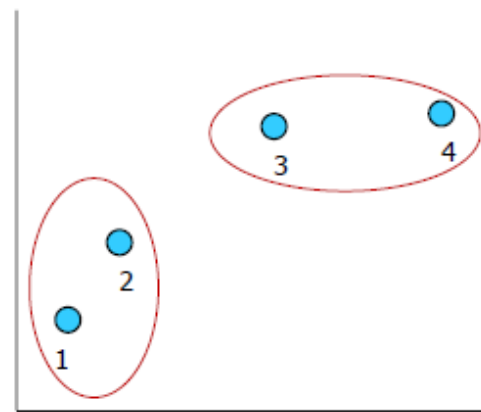
A dissimilaridade entre dois clusters (possivelmente **singletons**) é representada como a altura do nó interno mais baixo compartilhado



Dendrograma



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

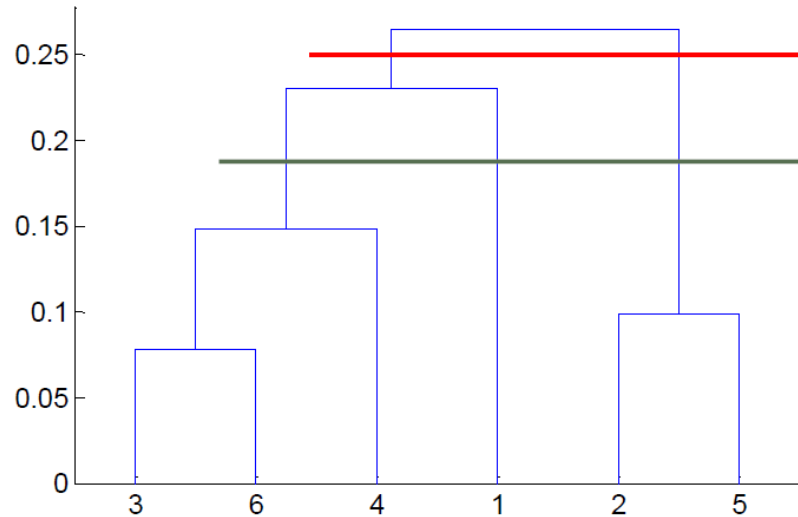


Dendrograma -> Grupos



Partições são obtidas via **cortes** no dendrograma

- cortes horizontais
- no. de grupos da partição = no. de interseções



$$G_1 = \{ (x1, x3, x4, x6), (x2, x5) \}$$

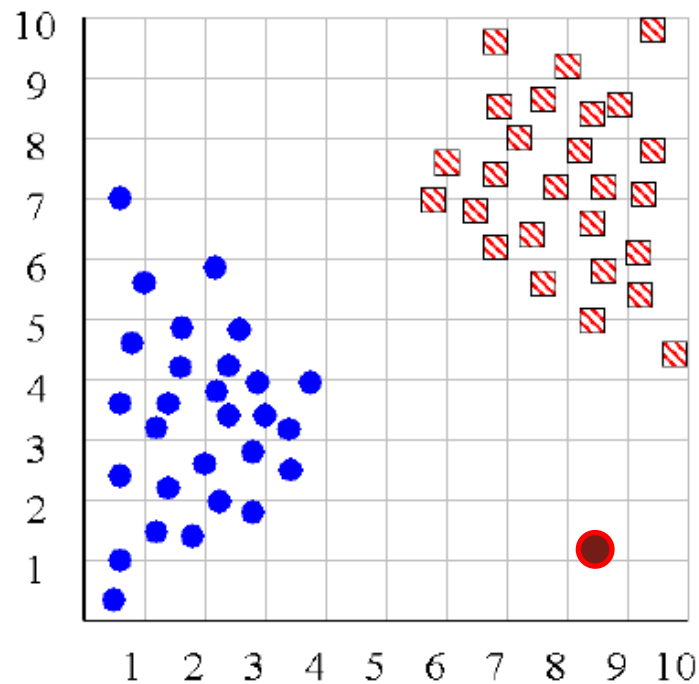
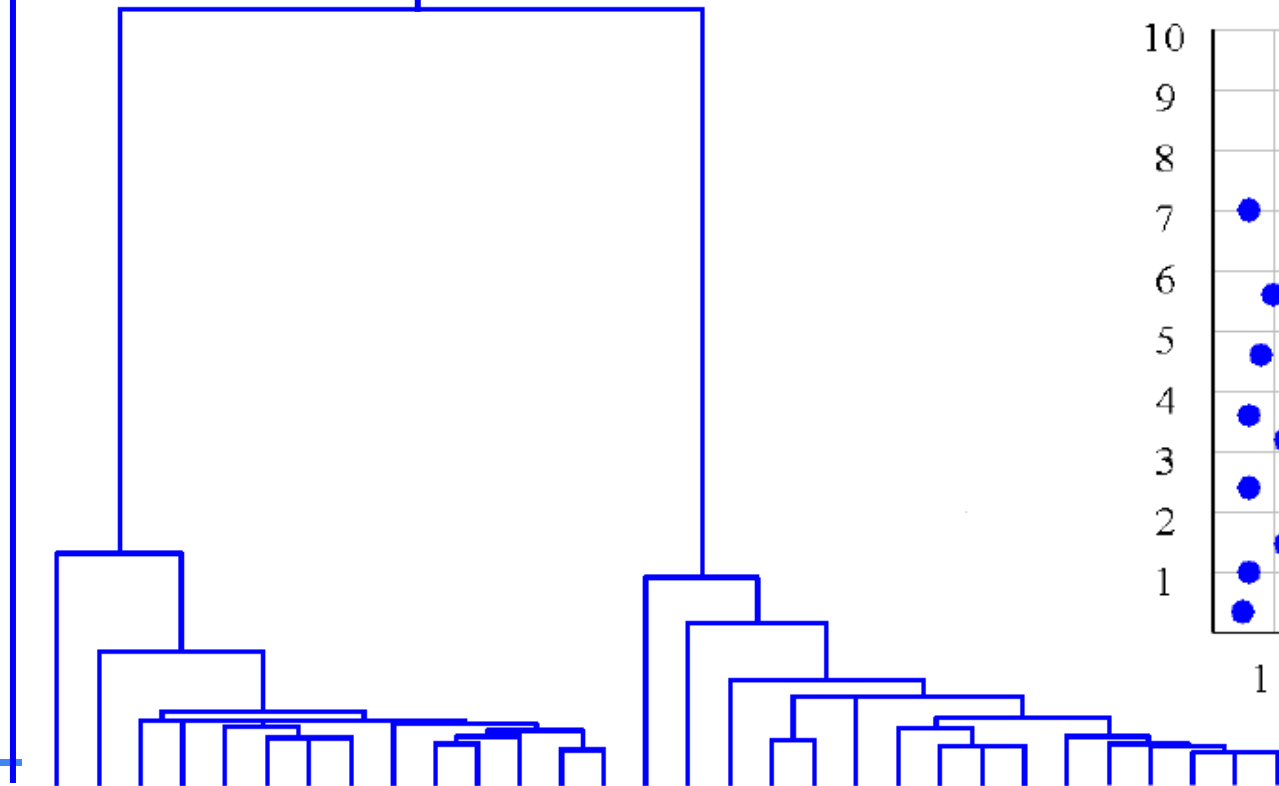
$$G_2 = \{ (x1), (x3, x4, x6), (x2, x5) \}$$

The plot shows the distribution of the number of nodes in the largest component of a network. The x-axis represents the number of nodes (1 to 10), and the y-axis represents the frequency (1 to 10). The plot shows a single bar at x=1 with a height of 10, and a single dot at (10, 7).



Dendrograma -> Outlier

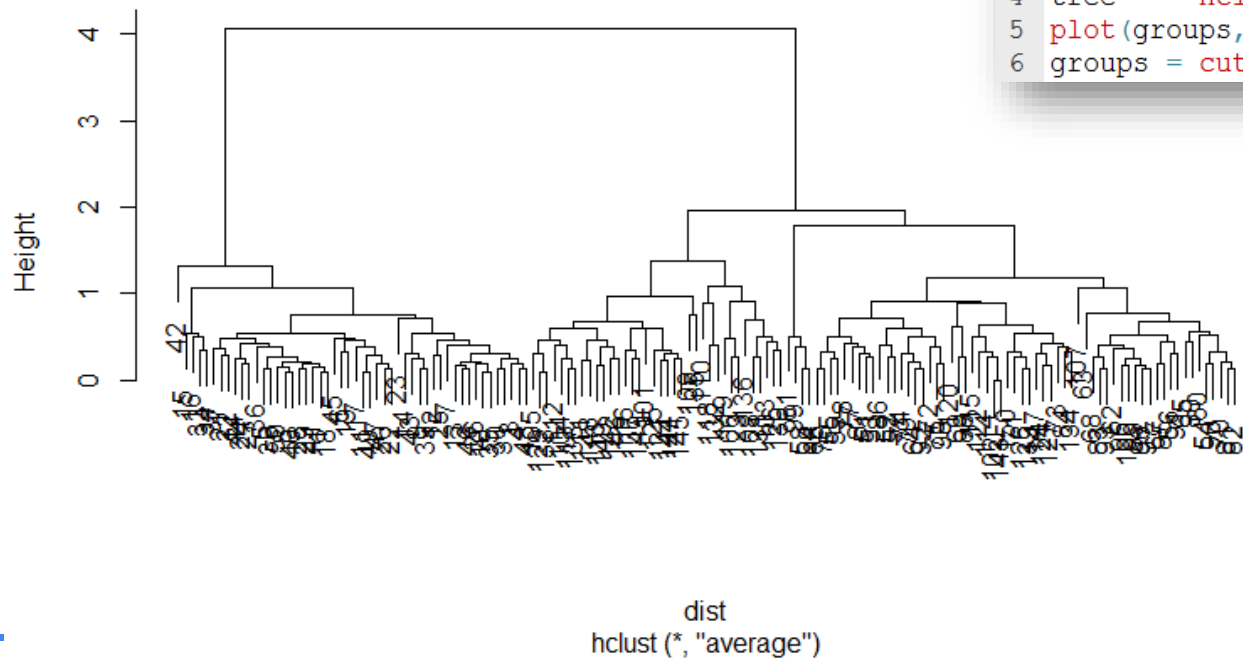
Pode-se examinar o dendrograma para tentar detectar a presença de outliers.



Hierárquico



Cluster Dendrogram



```
1 #Carrega os dados
2 data(iris)
3 dist = dist(iris[1:4])
4 tree = hclust(dist, method = "average")
5 plot(groups, col=unclass(iris$Species))
6 groups = cutree(tree, k=3)
```



Questions and Feedback



[Thank you!](#)

Obrigado !



Vinicius Fernandes Caridá
vfcarida@gmail.com



@Vinicius Caridá



@Vinicius Caridá
@machine learning Brasil



@vfcarida

MBA⁺

Copyright © 2018 **Prof. Vinicius Fernandes Caridá**
 Todos direitos reservados. Reprodução ou divulgação
 total ou parcial deste documento é expressamente
 proibido sem o consentimento formal, por escrito, do
 Professor (autor).