

# CLUSTERING

---

Prof. Nielsen Rechia

[nielsen.machado@uniritter.edu.br](mailto:nielsen.machado@uniritter.edu.br)

# Clustering

2

Paradigmas	Supervisionado		Não-supervisionado	
Tarefas	Classificação		Análise associativa	
	Regressão		Agrupamento ( <i>clustering</i> )	
	Outros		Redução de dimensionalidade	
			Outros	

7 tarefas comuns de aprendizado de máquina:  
<http://vitalflux.com/7-common-machine-learning-tasks-related-methods/>

# Clustering

**No aprendizado não supervisionado normalmente não temos a informação da classe das instâncias de treinamento.**

Agrupamento é uma tarefa que visa agrupar um conjunto de objetos (instâncias) em diferentes classes (grupos) de objetos de acordo com sua similaridade.

# Clustering

## Problemas:

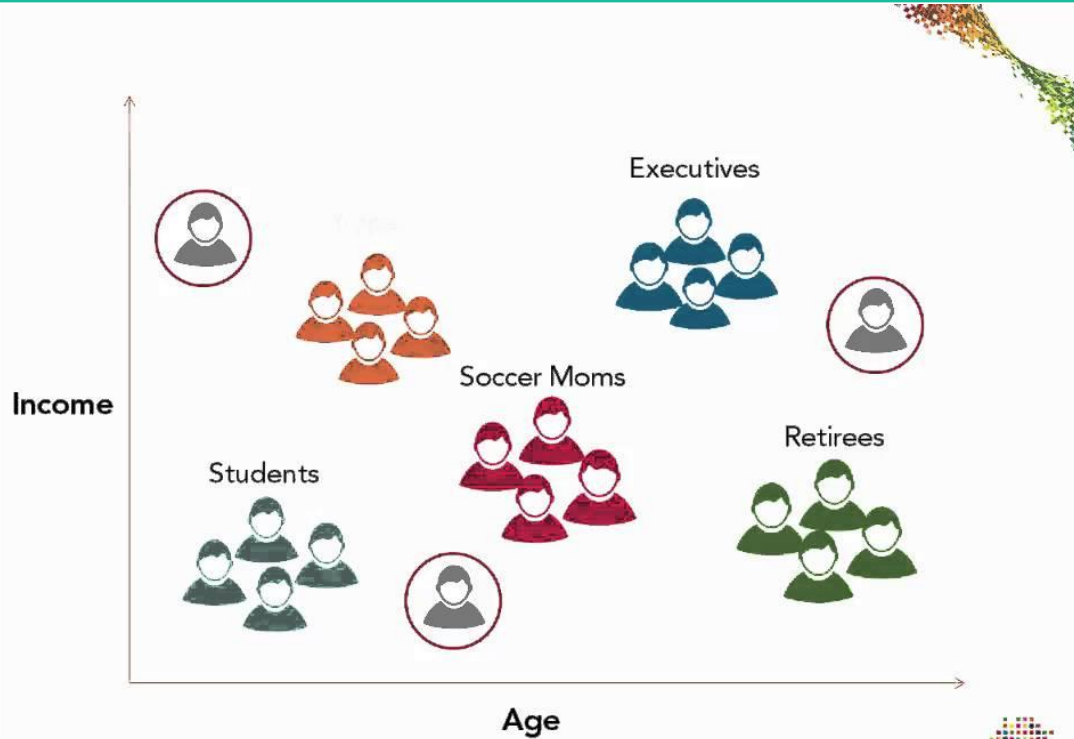
Os consumidores possuem perfis similares?

É possível agrupar tipos de câncer pelo comportamento que apresentam?

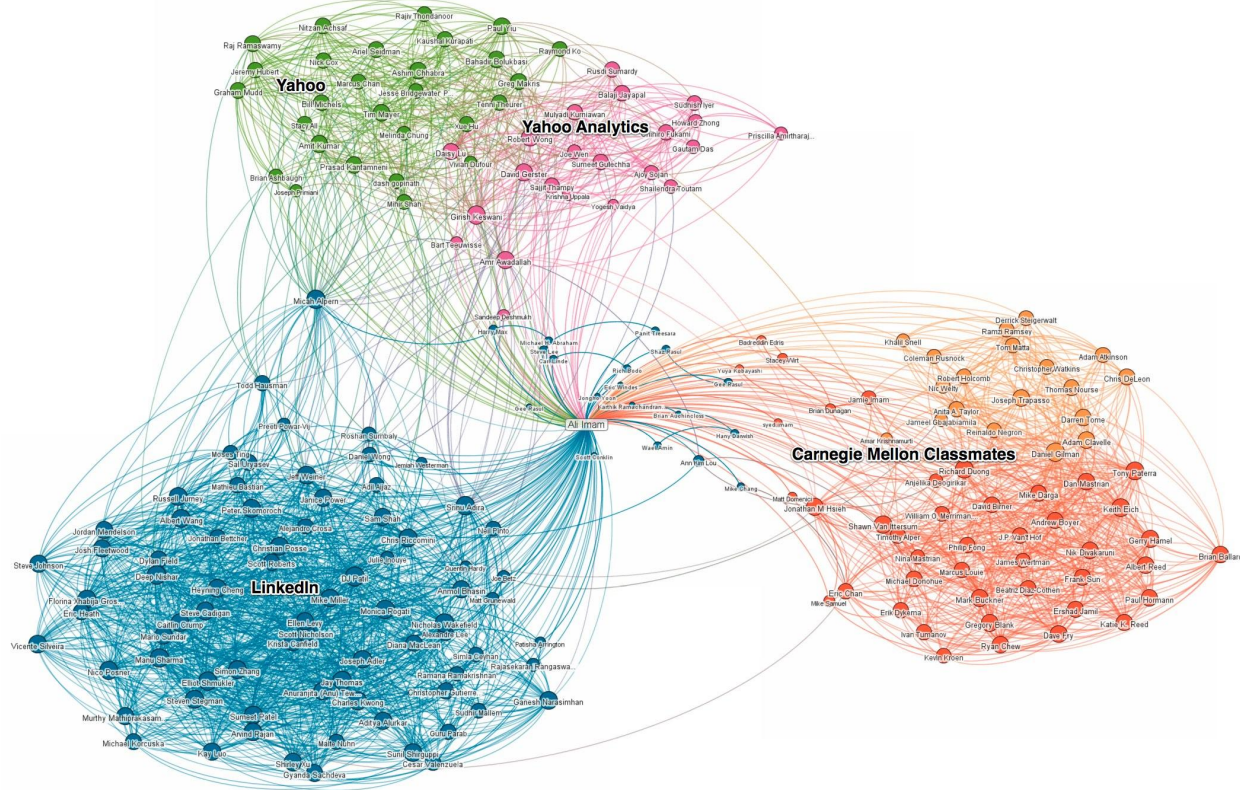
Quais as regiões onde determinado tipo de crime acontece mais frequentemente?

# Clustering

5



# Clustering



# Clustering

## Existem tipos de algoritmos de agrupamento

Algoritmos particionais:

K-Means

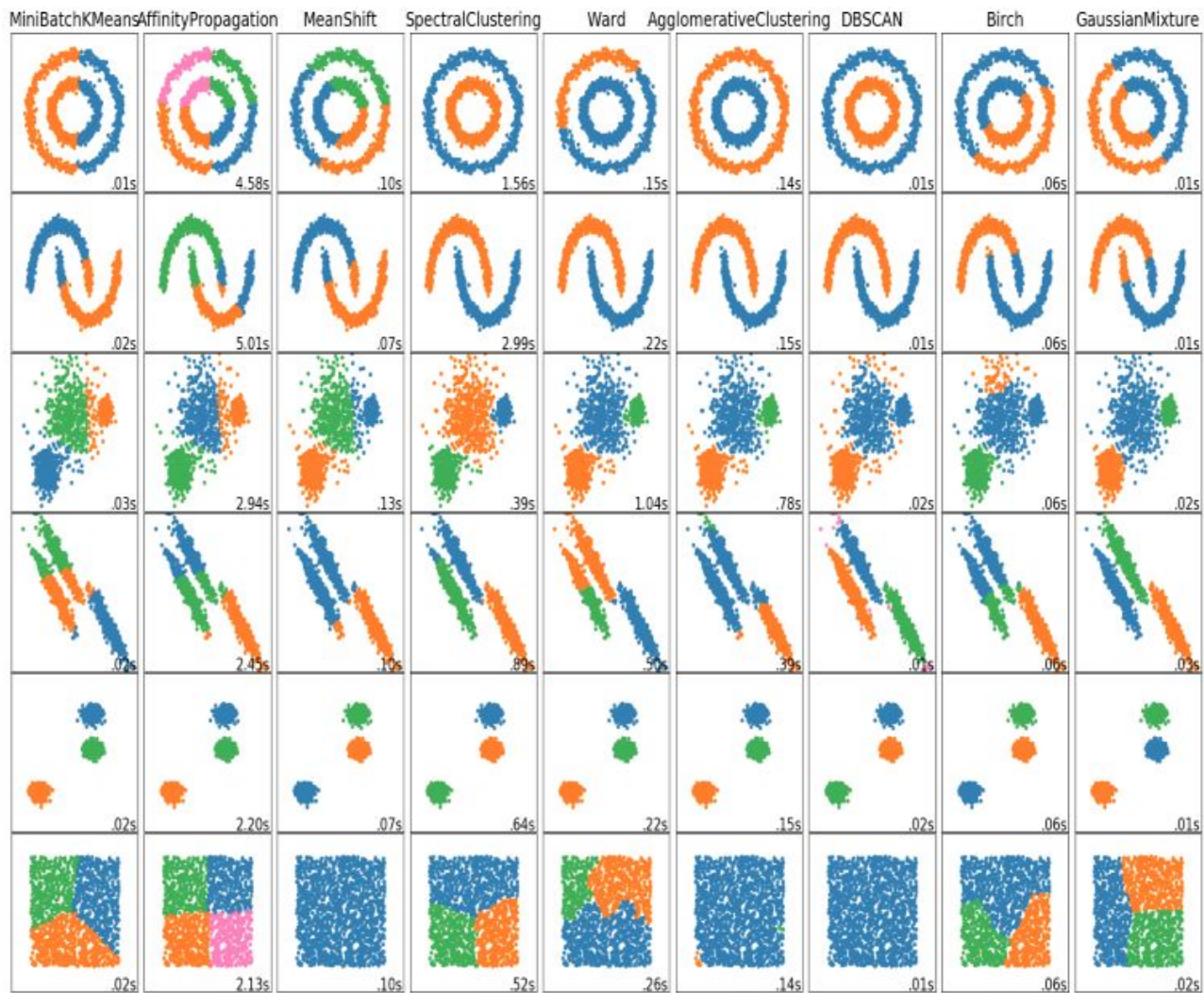
Mais popular, rápido e com bom desempenho

Algoritmos de densidade:

DBSCAN

Algoritmos hierárquicos:

Complete Linkage, Single Linkage, WARD, Average (UPGMA), PAM e ROCK





# Clustering

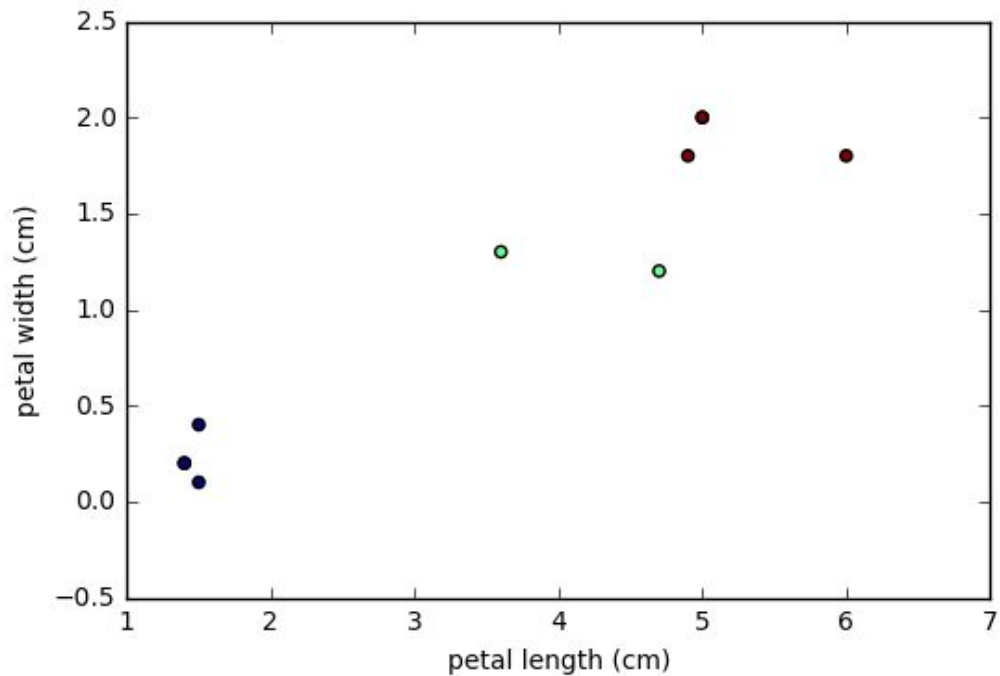
Dado parte do conjunto de dados iris, como agrupar as instâncias similares?

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	class
9	4.9	3.1	1.5	0.1	0.0
125	7.2	3.2	6.0	1.8	2.0
15	5.7	4.4	1.5	0.4	0.0
64	5.6	2.9	3.6	1.3	1.0
113	5.7	2.5	5.0	2.0	2.0
123	6.3	2.7	4.9	1.8	2.0
113	5.7	2.5	5.0	2.0	2.0
8	4.4	2.9	1.4	0.2	0.0
73	6.1	2.8	4.7	1.2	1.0
0	5.1	3.5	1.4	0.2	0.0

# Clustering

10

Dados originais no plano.



# Clustering - K-Means

11

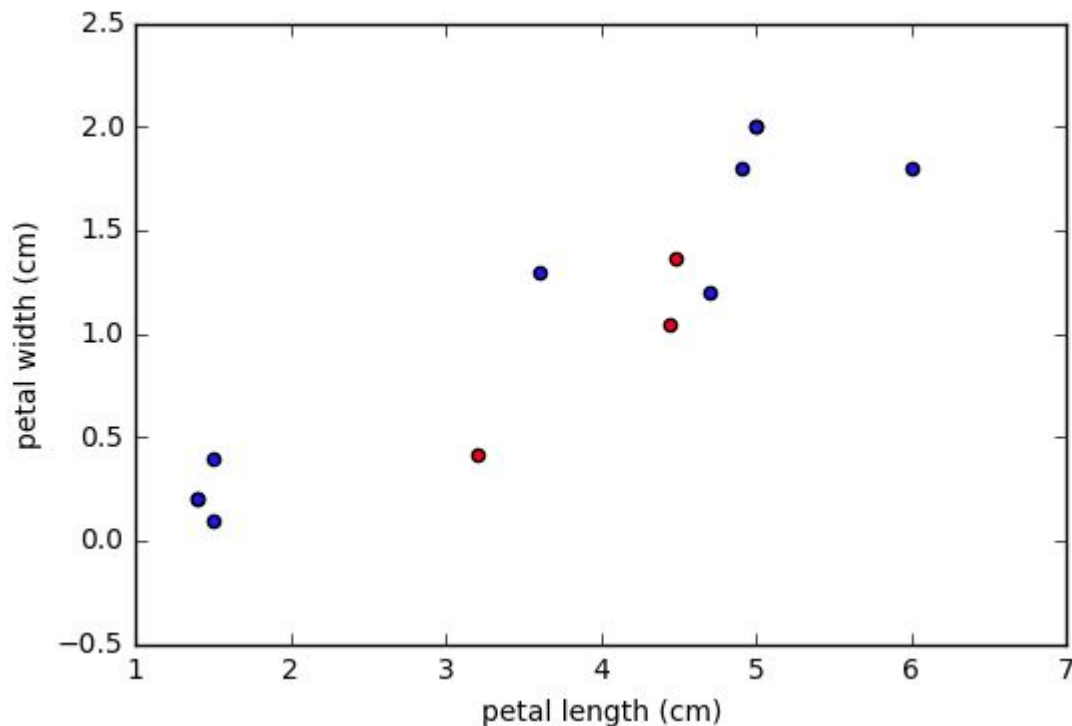
Escolhemos aleatoriamente K centróides para os clusters (grupos).

Intervalos de valores:

	petal length (cm)	petal width (cm)
min	1.4	0.1
max	6.0	2.0

centróides:

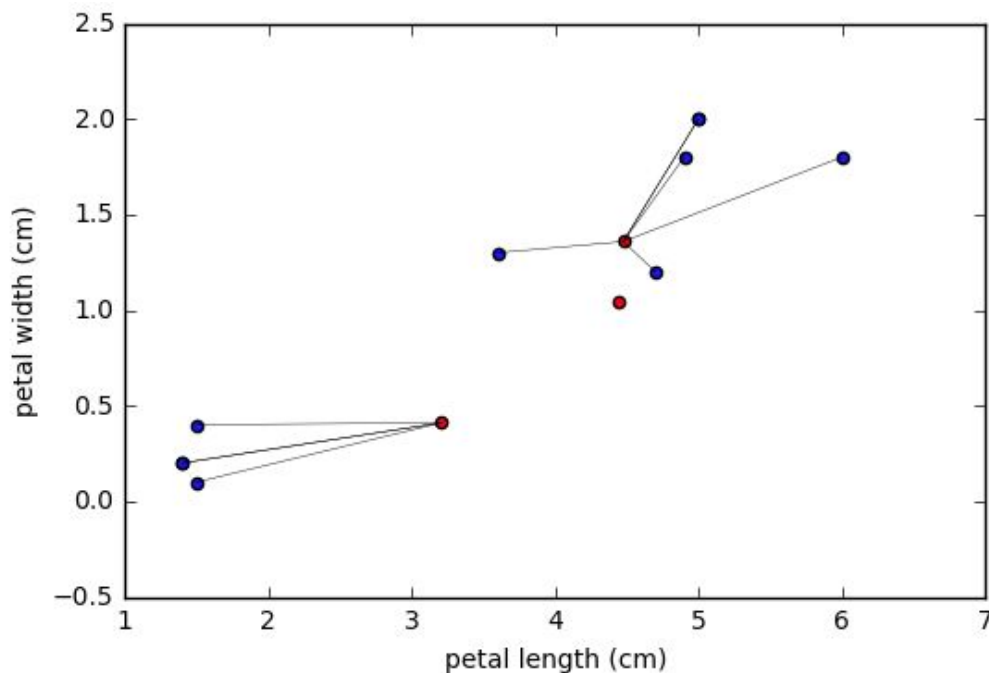
	petal length (cm)	petal width (cm)
0	4.448337	1.039741
1	4.474348	1.358150
2	3.209553	0.410570



# Clustering - K-Means

12

Atribuímos cada uma das instâncias ao centróide mais próximo.



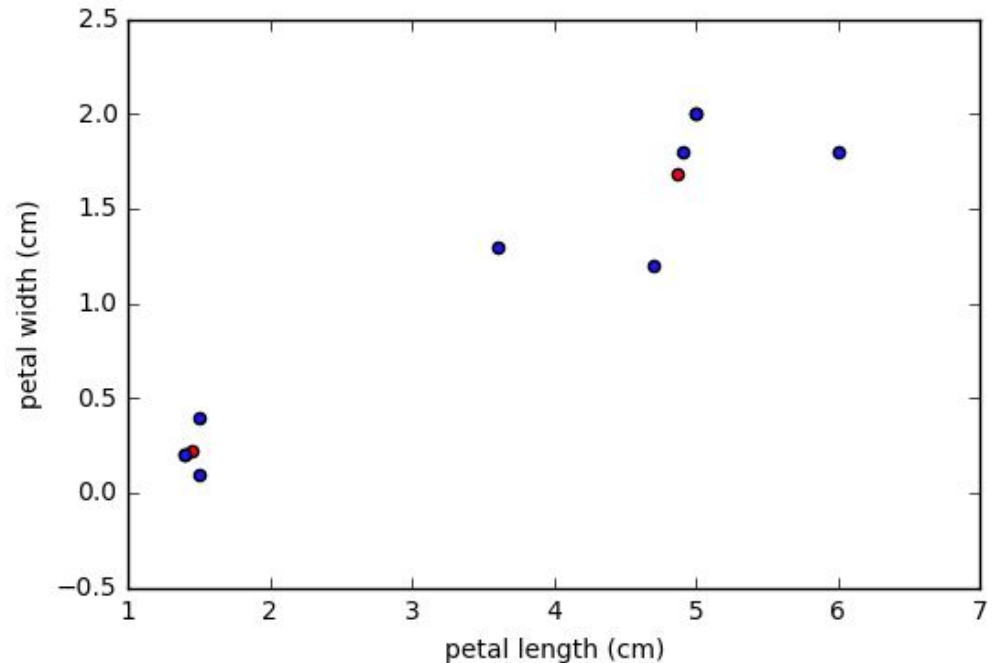
# Clustering - K-Means

13

Movemos cada centróides para a média dos objetos do cluster correspondente

novos centróides:

	petal length (cm)	petal width (cm)
0	NaN	NaN
1	4.866667	1.683333
2	1.450000	0.225000



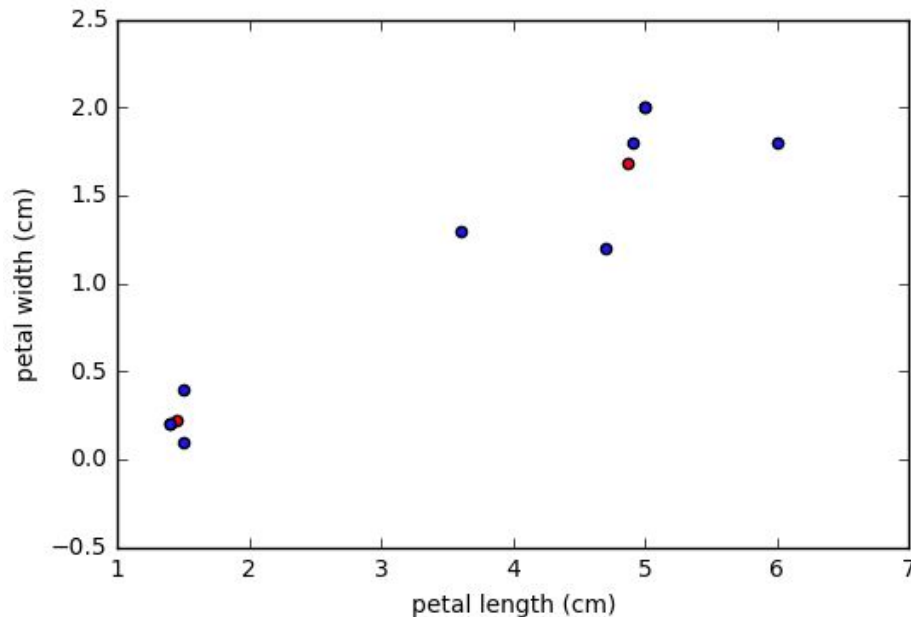
# Clustering - K-Means

14

Iteramos entre os passos anteriores até que algum critério de convergência seja satisfeito

Número máximo de iterações seja atingido

Limiar mínimo de mudanças no centróides



# Clustering - K-Means

15

Minimizar a seguinte função de custo:

$$J = \sum_{c=1}^k \sum_{\mathbf{x}_j \in C_c} \left\| \mathbf{x}_j - \bar{\mathbf{x}}_c \right\|^2$$

j-ésima instância do c-ésimo grupo

Distância euclidiana ao quadrado

Para cada grupo  $C_c$

Para cada instância no grupo  $C_c$

Posição do centróide do c-ésimo grupo

# Clustering - K-Means

16

**K-Means pode sofrer problemas ao lidar com grupos de diferentes**

- Tamanhos

- Densidades

- Formas (principalmente não-esféricas)

- É sensível a outliers

Existem variações como **K-medianas**



# Hierarchical Clustering

17

## Algoritmos hierárquicos

Ex: single linkage, ward, complete linkage

Permitem verificar a relação que os dados possuem

Relação é visual em uma estrutura chamada dendrograma

### Existem duas formas de agrupamento hierárquico:

**Aglomerativo:** Começa com N grupos de uma instância, e combina tais grupos até chegar a um único grupo

**Divisivo:** Começa com apenas um grupo, e separa sucessivamente até terminar em grupos de uma instância

Na prática, algoritmos aglomerativos são mais populares

# Agglomerative Clustering

18

## Dendrograma

É uma ferramenta útil para sumarizar medidas de (dis)similaridade e mostrar as hierarquias de partições resultantes de algoritmos hierárquicos aglomerativos

Pode-se examinar o dendrograma inclusive para estimar o número correto de grupos.

Partições são obtidas via cortes no dendrograma

Cortes horizontais

No de grupos da partição = Número de intersecções

# Agglomerative Clustering

19

## Algoritmos hierárquicos

Ex: single linkage

$$d(u, v) = \min(\text{dist}(u[i], v[j]))$$

Distância entre os  
clusters  $u, v$

$i$ -ésima instância do  
cluster  $u$

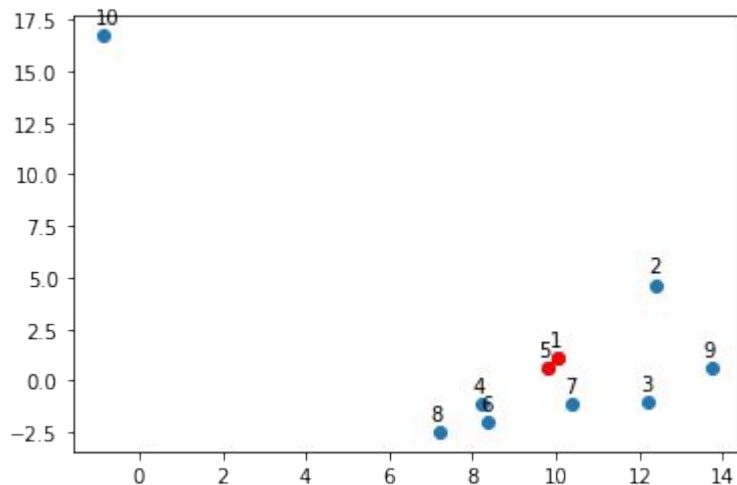
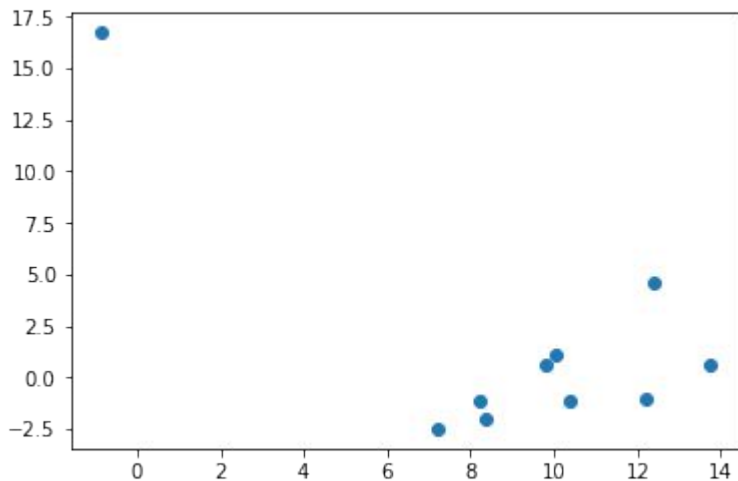
$j$ -ésima instância do  
cluster  $v$

# Agglomerative Clustering

20

## Single Linkage

Encontrar a menor distância entre os objetos do conjunto de dados

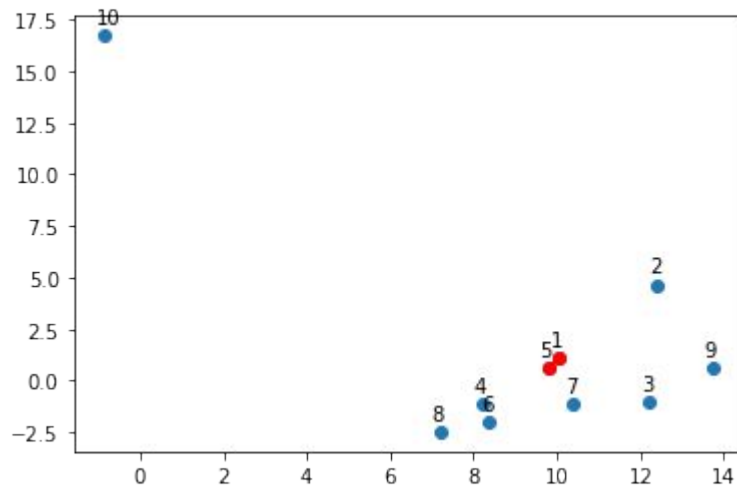


# Agglomerative Clustering

21

## Dendrograma

Conecte os dois grupos (instâncias) da distância encontrada

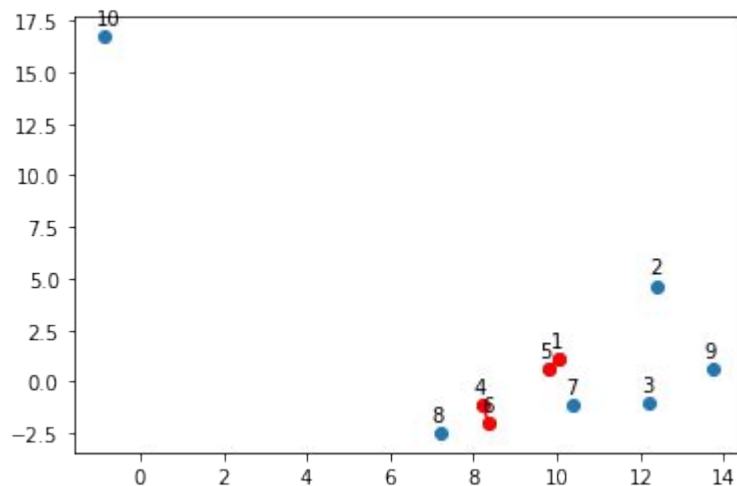
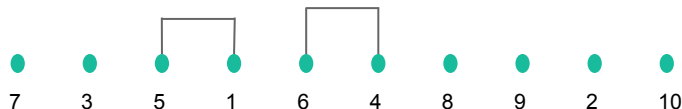


# Agglomerative Clustering

22

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

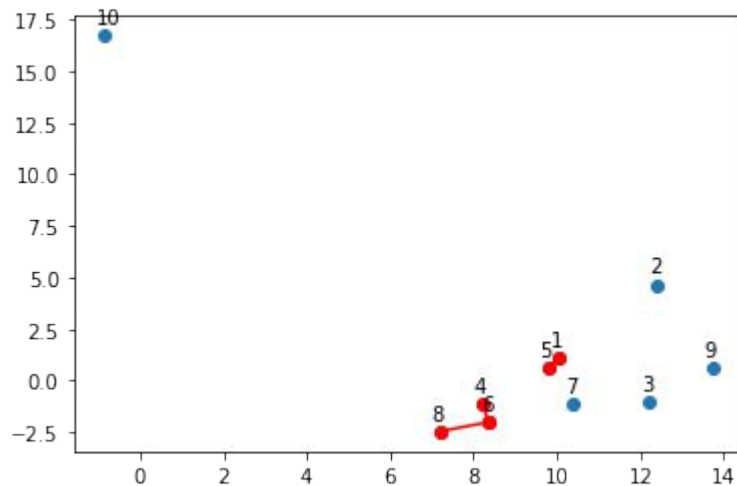
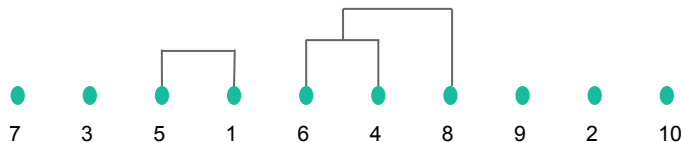


# Agglomerative Clustering

23

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

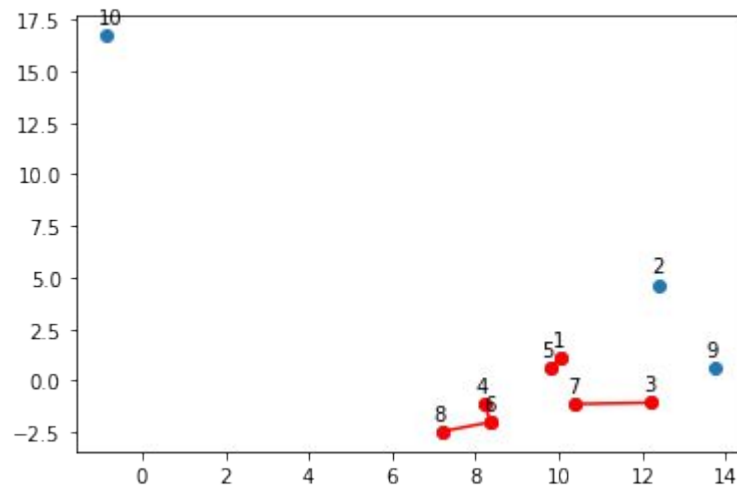
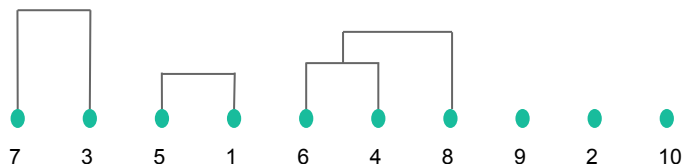


# Agglomerative Clustering

24

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas



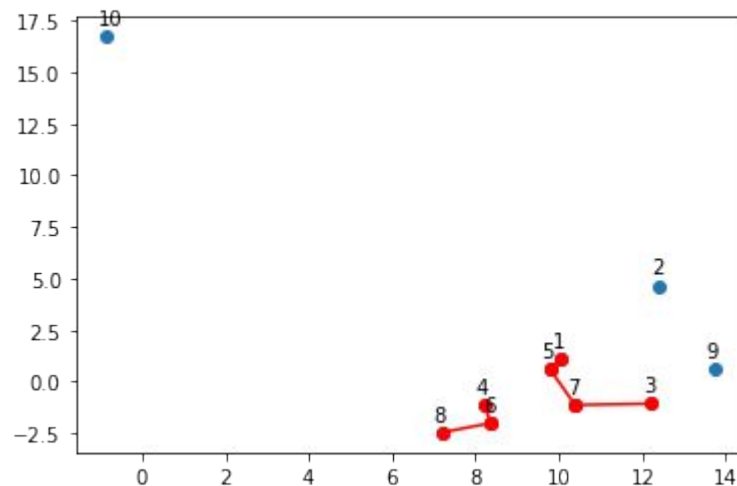
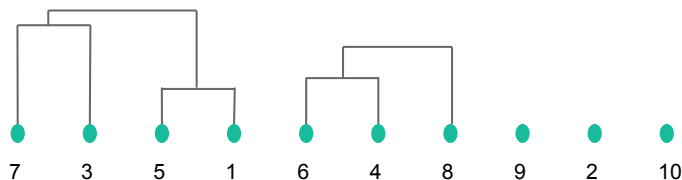


# Agglomerative Clustering

25

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

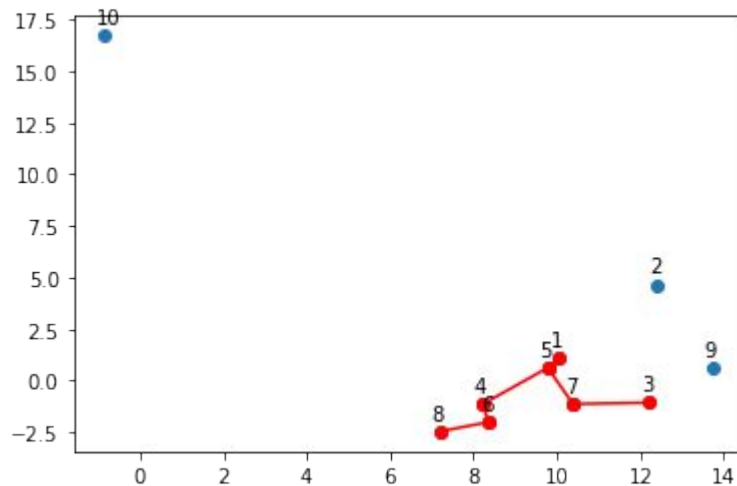
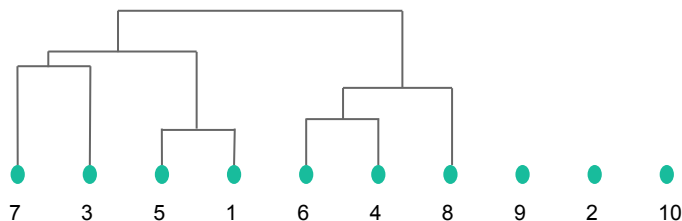


# Agglomerative Clustering

26

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

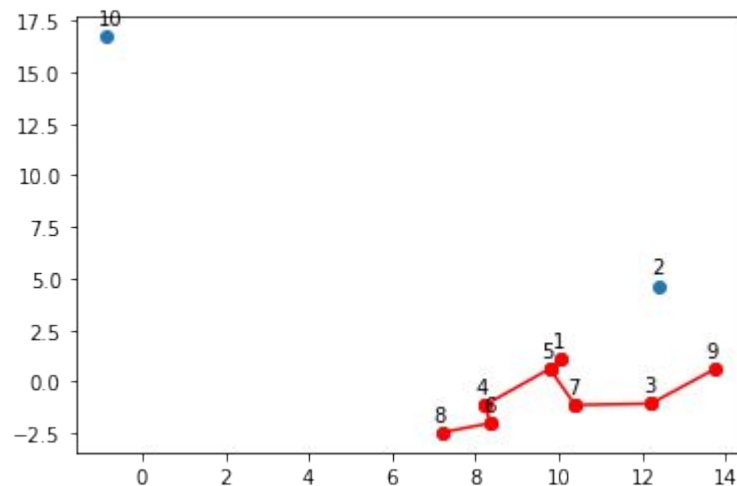
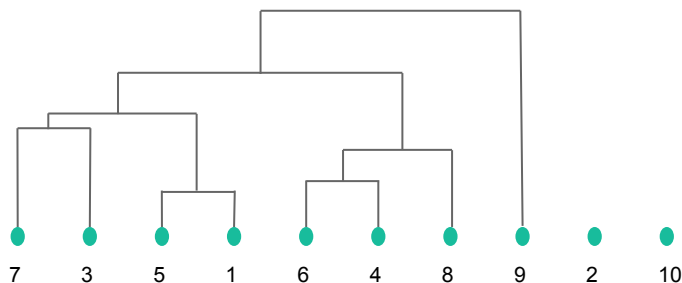


# Agglomerative Clustering

27

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

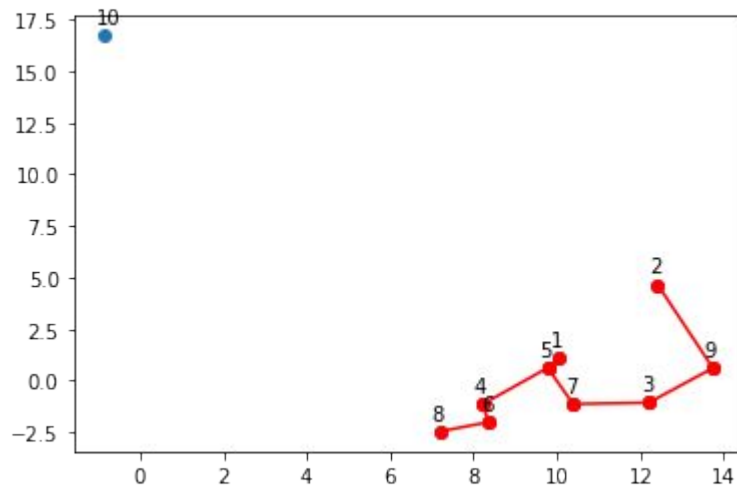
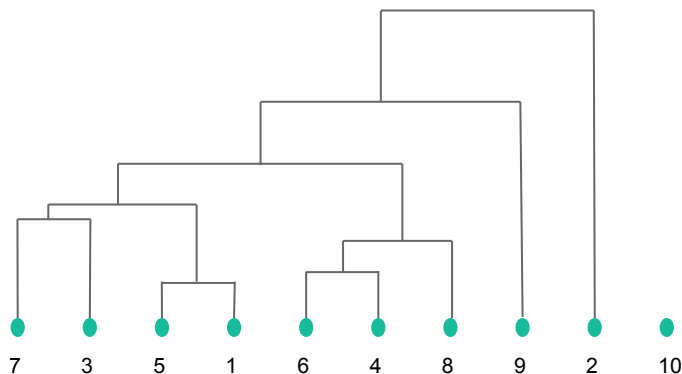


# Agglomerative Clustering

28

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas

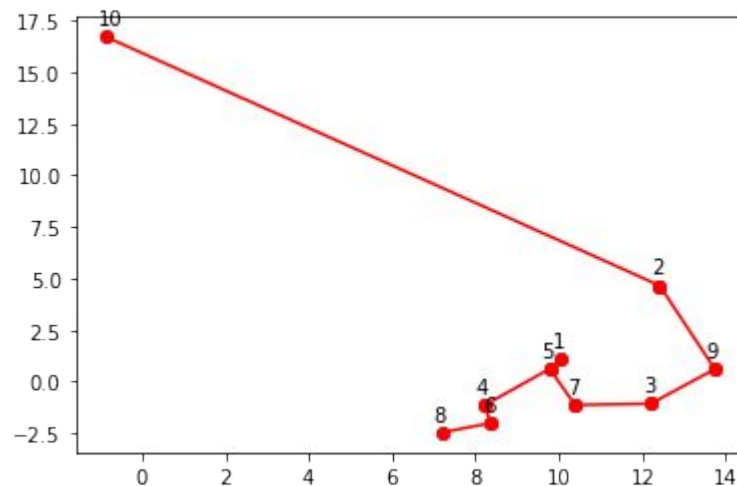
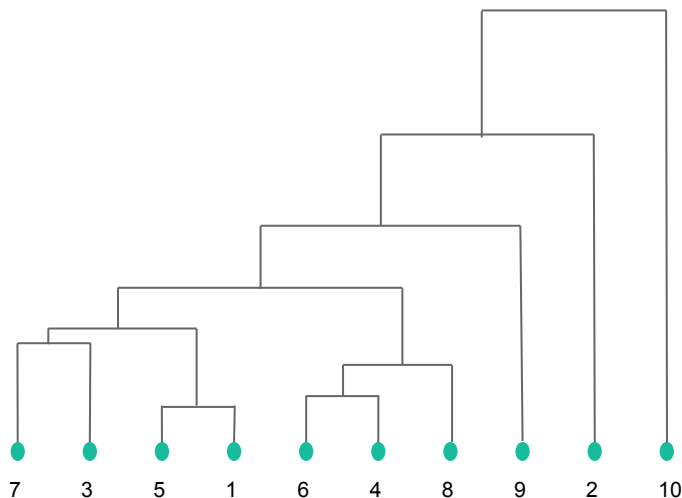


# Agglomerative Clustering

29

## Dendrograma

Repita o processo até que todas as instâncias estejam conectadas



# Agglomerative Clustering

30

**Para dividir os dados em grupos:**

Realize um corte no dendrograma onde a distância entre os dois grupos for a maior

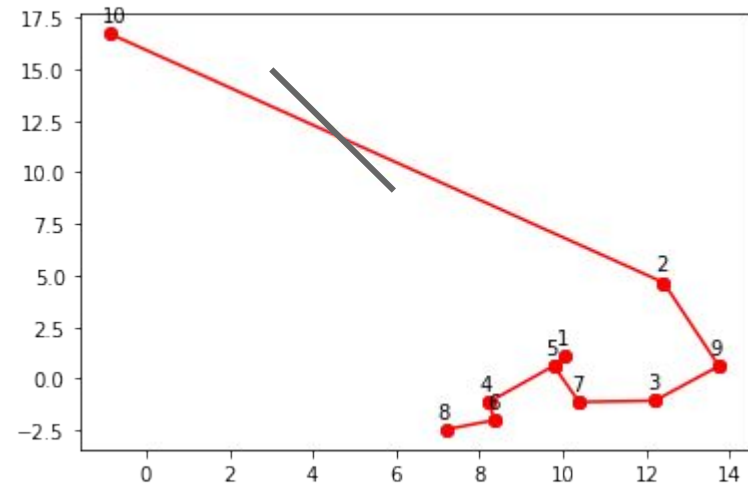
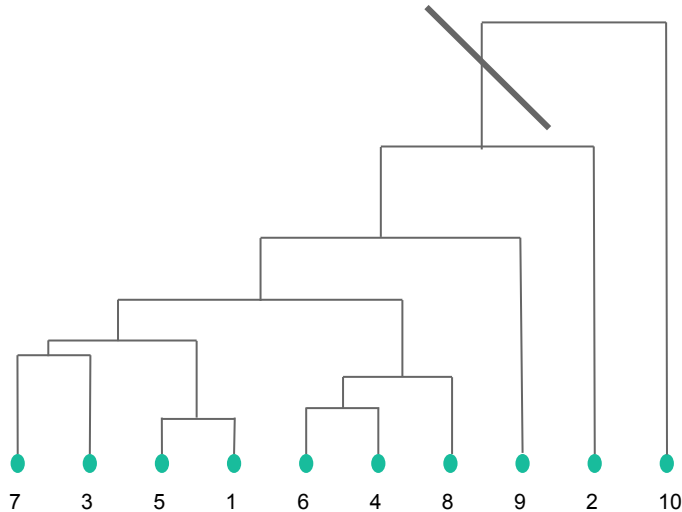
**Uma boa regra para realizar agrupamentos é que os dados devem possuir bastante similaridade intra-grupos e pouca similaridade inter-grupos**

Caso este critério não seja atingido, é possível continuar realizando cortes

Geralmente utilizam-se medidas de avaliação para verificar este critério

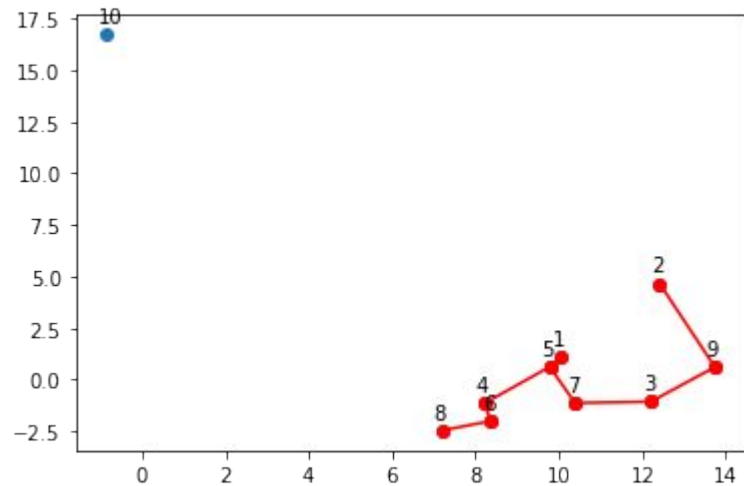
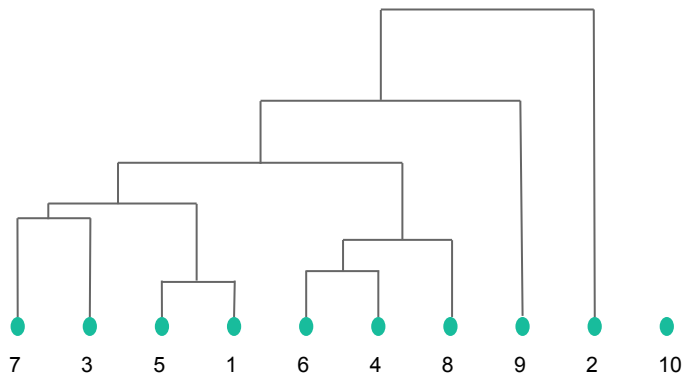
# Agglomerative Clustering

31



# Agglomerative Clustering

32





# Agglomerative Clustering

33

Avalia quantitativamente uma partição/ hierarquia de dados a partir de alguma medida. São divididos em 3 tipos:

**Externos:** avaliam o grau de correspondência entre a estrutura encontrada e uma solução esperada ou conhecida (golden truth).  
Ex: Rand Index, Jaccard, Fowlkes-Mallows

**Internos:** avaliam o grau de compatibilidade entre a estrutura encontrada e os dados (e apenas os dados). Ex: SSE

**Relativos:** avaliam dentre duas ou mais estruturas, qual a melhor (segundo algum aspecto). Tipicamente são índices internos capazes de quantificar a qualidade dos agrupamentos. Ex: Davis-Bouldin, Silhueta, Dunn, Gap Statistics, Calinski-Harabasz

# Clustering - Avaliação

34

## Medidas de avaliação de clustering

Silhueta (Silhouette Width Criterion):

Para cada instância,  
calcula o índice  $s(i)$

$$SWC = \frac{1}{N} \sum_{i=1}^N s(i)$$

Dissimilaridade média  
da  $i$ -ésima instância  
ao segundo cluster  
mais próximo

$$s(i) = \frac{b(i) - a(i)}{\max\{a(i), b(i)\}}$$

Dissimilaridade  
média da  $i$ -ésima  
instância ao cluster  
mais próximo

A swc mede quão compactos e bem-separados os dados estão

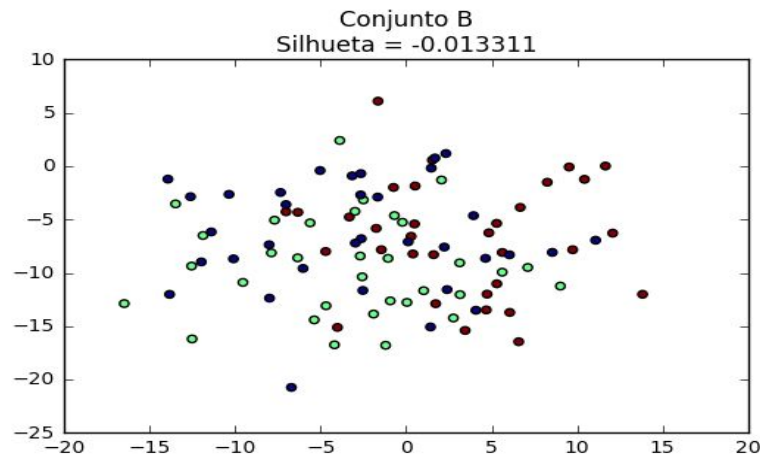
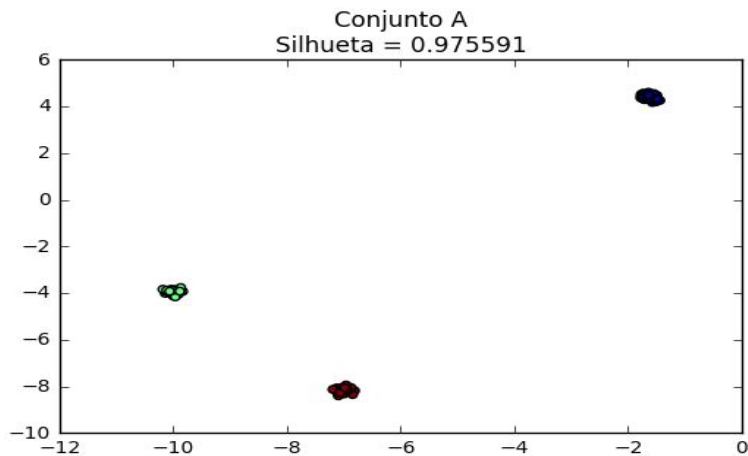
# Clustering - Avaliação

A silhueta varia de -1 até 1:

-1: grupos esparsos e misturados

1: grupos bem separados e compactos

Entre 0 e 0.5 indica que os dados não são bem separados



# Clustering - Avaliação

36

Mais métricas para avaliar agrupamentos:

<http://scikit-learn.org/stable/modules/classes.html#clustering-metrics>

K-Means - sklearn:

```
from sklearn.cluster import KMeans
cl = KMeans(n_clusters=3, init='k-means++', n_init=10, max_iter=300)
cl = cl.fit(df)

labels = cl.labels_
centroids = cl.cluster_centers_
```

Hierarquicos - scipy:

```
from scipy.cluster.hierarchy import dendrogram, linkage
from scipy.cluster.hierarchy import fcluster
Z = linkage(dataset, 'single')
```

# Exercício prático

37

**Implemente o algoritmo k-means e single linkage para o dataset iris.**

Remova a classe do conjunto de dados

Execute com diversos números de grupo (ex, de 2 a 10)

Execute apenas uma vez para cada k

Verifique os agrupamentos utilizando a silhueta

Qual é o melhor número de grupos?

Qual é o melhor algoritmo? PQ?

```
In [77]: from sklearn.datasets import load_iris
import numpy as np
import pandas as pd

data = load_iris()
df = pd.DataFrame(
    data['data'],
    columns=data['feature_names']
)
```

# Exercício para entregar

38

**Agrupe os dados do censo americano de 2005 com o K-means e algum algoritmo hierárquico**

Não considerar a classe do conjunto de dados

Somente um conjunto de dados

Dados devem ser pré-processados

Usar vários números de grupos (ex: 2 a 20)

Medir a qualidade dos grupos com a silhueta (podem usar outros também)

Qual é o melhor número de grupos? Qual é o melhor algoritmo?  
Motivos?

Entregar via BB o código python ou jupyter e o passo a passo realizado (ex, motivos do pre-processamento etc.) em um zip.

# Conclusão

39

## Leitura recomendada:

Capítulo 8 e 9 de Introduction to Data Mining

